# BMJ Open

## Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm Assessing Dementia Risk in the Community

SCHOLARONE™
Manuscripts

1
2
3   1   **Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm**
4
5   2   **Assessing Dementia Risk in the Community**
6
7
8   3
9

10
11  4   Stacey Fisher, MSc[1,2,3] stacey.fisher@uottawa.ca
12
13  5   Amy Hsu, PhD[1,2] ahsu@toh.ca
14
15  6   Nassim Mojaverian, MSc[2] namojaverian@ohri.ca
16
17  7   Monica Taljaard, PhD[1,3] mtaljaard@ohri.ca
18
19
20  8   Greg Huyer, MSc[1,3] ghuye047@uottawa.ca
21
22  9   Doug Manuel, MD[1,2,3,4,5] dmanuel@ohri.ca
23
24
25  10  Peter Tanuseputro, MD[1,2,5,6,7] ptanuseputro@ohri.ca
26
27  11
28
29
30  12  [1] Ottawa Hospital Research Institute, Ottawa, Ontario, Canada
31
32  13  [2] Institute for Clinical Evaluative Sciences, Ottawa, Ontario, Canada
33
34  14  [3] Department of Epidemiology and Community Medicine, University of Ottawa, Ottawa,
35
36  15  Ontario, Canada
37
38
39  16  [4] Statistics Canada, Ottawa, Ontario, Canada
40
41  17  [5] Department of Family Medicine, University of Ottawa, Ottawa, Ontario, Canada
42
43
44  18  [6] Department of Medicine, University of Ottawa, Ottawa, Ontario, Canada
45
46  19  [7] Bruyère Research Institute, Ottawa, Ontario, Canada
47
48  20
49
50
51  21  Corresponding author:
52
53  22  Dr. Peter Tanuseputro
54
55
56  23  1053 Carling Ave.
57
58
59
60

24    Box 693

25    Ottawa, ON K1Y 4E9, Canada

26    ptanuseputro@ohri.ca

27    Phone: 613-798-5555

28

29    **Word Count:** 3 723

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

2

**ABSTRACT**

**Introduction:** The burden of disease from dementia is a growing global concern as incidence increases exponentially with age and average life expectancy has been increasing around the world. Planning for an aging population requires reliable projections of future dementia prevalence and resource requirements, however, existing population projections are simple and have poor predictive accuracy. The Dementia Population Risk Tool (DemPoRT) will predict incidence of dementia in the population setting using multivariable modeling techniques.

**Methods and Analysis:** The derivation cohort will consist of elderly Ontario respondents of the Canadian Community Health Survey (CCHS) (2001, 2003, 2005, 2007; 18 764 males and 25 288 females). Pre-specified predictors include sociodemographic, general health, behavioral, functional and health condition variables. Incident dementia will be identified through individual linkage of survey respondents to population-level administrative health care databases (1 797 and 3 281 events, and 117 795 and 166 573 person-years of follow-up, for males and females, respectively until March 31, 2014). Using time of first dementia capture as the primary outcome and death as a competing risk, sex-specific proportional hazards regression models will be estimated. The 2008/2009 CCHS survey will be used for validation (approximately 4 600 males and 6 300 females). Overall calibration and discrimination will be assessed as well as calibration within predefined subgroups of importance to clinicians and policy makers.

**Ethics and Dissemination:** This study has been approved by the Ottawa Health Science Network Research Ethics Board. DemPoRT results will be submitted for publication in peer-review journals and presented at scientific meetings. The algorithm will be assessable online for both population and individual uses.

**Trial Registration Number:** ClinicalTrials.gov NCT03155815.

70 **STRENGTHS AND LIMITATIONS**

71 - The Dementia Population Risk Tool (DemPoRT) will be developed and validated using

72 predictors from large population-based community health surveys that are individually

73 linked to routinely-collected health administration data in Ontario. To our knowledge,

74 DemPoRT will be the first population-based algorithm for predicting and projecting

75 dementia incidence.

76 - DemPoRT will produce improved estimates of future dementia burden, will assess the

77 contribution of specific risk factors to the population risk, and will identify population

78 subgroups at high risk of developing dementia. This information will be used by

79 policymakers to prepare for and reduce dementia impact.

80 - The analysis plan and predictors have been fully pre-specified to limit the risk of over-

81 fitting and improve the quality of predictions.

82 - Detailed cognitive testing to ascertain dementia diagnoses is preferable over the use of

83 administrative data, however this is not available or feasible at the population level.

84 - Although a rigorous approach to model development will be used, further validation will

85 be needed to assess generalizability, and calibration will be required for application in

86 other jurisdictions.

87

88

89

90

91

92

93      **INTRODUCTION**

94      The burden of disease from dementia is a growing global concern as incidence increases

95      exponentially with age and average life expectancy has been increasing around the world[1,2].

96      Planning for an aging population requires reliable projections of future dementia prevalence and

97      its implications on resource requirements. Existing population projections for dementia,

98      however, are overly simplistic and likely inaccurate[3].

99

100     **Limitations of Current Dementia Projection Methodology**

101     Almost all existing dementia projections have used extrapolation and macrosimulation methods,

102     which are simplistic and make assumptions that may not hold true into the future[3]. Most

103     extrapolations simply apply current age- and sex-specific prevalence of dementia to future

104     population projections. Macrosimulations use estimates of dementia incidence and mortality,

105     stratified by age and sex, to simulate disease prevalence as the population ages[1,4–6]. Projections

106     from extrapolations incorrectly assume that the risk of mortality among those with and without

107     dementia are equivalent[7,8], and both methods assume that the age and sex-specific prevalence of

108     dementia risk factors will not change with time. The assumption of stable risk factor prevalence

109     is widely thought to be the major source of error in existing dementia projections[3,9–11].

110

111     Up to 50% of dementia cases may be attributable to physical inactivity, obesity, diabetes,

112     hypertension, low educational achievement and depression[9,12]. Changing trends of these risk

113     factors over time has the potential to have a strong impact on dementia prevalence. For example,

114     the population prevalence of diabetes and obesity in Canada has been projected to increase,

115     while smoking, hypertension and dyslipidemia have been projected to decline[13]. Consideration of

116  risk factor prevalence is therefore important to improve the accuracy of dementia projections,

117  and simple extrapolations and macrosimulations are often inadequate to incorporate changing

118  risk factors.

119

120  **Predictive Multivariable Modeling of Dementia Incidence**

121  Another method of dementia projection involves the development of population-based

122  **predictive risk algorithms** that examine the effect of risk factors on dementia incidence.

123  Population-based data that contain detailed exposure information, such as health surveys, are

124  linked at the individual-level to administrative data that capture dementia development. A

125  multivariable model of dementia incidence is derived, validated against external data, and

126  predictive performance is assessed. Counterfactual risk factor levels can be entered in to the

127  algorithm at the population level, or at individual level and summed, to simulate future disease

128  prevalence under different assumptions.

129

130  Incorporation of predictive risk algorithms in to microsimulation models such as Statistics

131  Canada's Population Health Models (POHEM) provides additional utility. POHEM dynamically

132  models individual life trajectories of a population representative of Canada including births,

133  deaths and migration, disease incidence and progression, and exposure to risk factors.

134  This facilitates detailed examination of the influence of changing risk factor prevalence on future

135  dementia prevalence and the potential influence of dementia prevention strategies to reduce the

136  population risk. In addition, these algorithms can be used to describe the risk of dementia in the

137  population, assess the contribution of specific risk factors to the population risk, and identify

138  high-risk groups.

139

140 The objective of this study is to develop and validate the Dementia Population Risk Tool

141 (DemPoRT) algorithm to predict dementia incidence in the population setting. This will be done

142 using multivariable modelling techniques, linking self-reported risk factors captured by a

143 population-based health survey in Canada with administrative databases across healthcare sectors

144 that capture healthcare diagnosed dementia. To our knowledge, the DemPoRT predictive model

145 will be the first population-based algorithm for predicting and projecting dementia incidence. It

146 will be able to estimate the future burden of dementia using techniques that consider changes in

147 risk factor prevalence and will identify modifiable risk factors that can be targeted by

148 individuals, clinicians and policy makers to reduce dementia incidence more effectively.

149

150 **METHODS AND ANALYSIS**

151 **Study Design**

152 Two DemPoRT models, one for males and females, will be derived and validated using

153 population-based data in Ontario, Canada, a multicultural province with 13.6 million residents.

154 Predictors will be obtained from the Canadian Community Health Surveys (CCHS), and

155 outcomes (i.e., diagnosis of dementia) will be obtained from routinely-collected health care data.

156

157 The derivation cohorts will consist of eligible respondents of the 2001, 2003, 2005 and 2007

158 CCHS (Cycles 1.1, 2.1, 3.1 and 4.1), while validation cohorts will consist of respondents to the

159 2008/2009 cycle. The CCHS is a national, cross-sectional survey developed by Statistics Canada

160 to collect information related to health and health care utilization of the Canadian population.

161 The survey has a multistage stratified cluster design that represents approximately 98% of the

162  Canadian population aged 12 years and over and attained an average response rate of 79% over

163  the study period. The CCHS is conducted through telephone and in-person interviews, and all

164  responses are self-reported. The details of survey methodology have been published elsewhere[14].

165  Survey respondents will be excluded if they are less than 55 years of age at survey

166  administration, self-reported a history of dementia, or are not eligible for Ontario's universal

167  health insurance. If a respondent was included in more than one CCHS cycle, only their earliest

168  survey response will be used.

169

170  **Outcome**

171  Survey respondents diagnosed with dementia will be identified through individual linkage to

172  several population-based administrative databases at the Institute for Clinical Evaluative

173  Sciences (ICES). Dementia case ascertainment is based on a validated definition: 1 hospital

174  record OR 3 physician claim records at least 30 days apart within a 2-year period OR a

175  dispensing record for a cholinesterase inhibitor from Ontario Drug Benefit (ODB). This

176  definition has a sensitivity of 79.3% and a specificity of 99.1% when validated against

177  emergency medical record (EMR) data[15]. Due to known underdiagnosis of dementia[16,17], we will

178  supplement this definition by adding survey respondents with dementia codes captured on home

179  care and long-term care assessments (dementia flag AND Cognitive Performance Scale [CPS]

180  score $\geq 2$) using the Resident Assessment Instrument-Home Care (RAI-HC) database and the

181  Continuing Care Reporting System (CCRS), respectively. We have found this addition adds

182  substantially (approximately 18%) to the number of dementia cases captured.

183

184  Survey respondents with dementia will be excluded if they meet the criteria for dementia within

185  2 years of survey administration (to remove potentially prevalent cases) or are younger than 65

186  years of age at the time of dementia diagnosis (to exclude early onset dementia which likely has

187  a different set of risk factors). Eligible survey respondents will be followed from the date of

188  survey administration or age 65, whichever came later, until the earliest date of: dementia

189  ascertainment, death (defined as competing risk), loss to follow-up (defined as loss of healthcare

190  eligibility) or end of study (March 31, 2014).

191

192  **Sample Size**

193  The male and female derivation cohorts consist of 18 764 and 25 288 respondents, and 117 795

194  and 166 573 person-years of follow-up, respectively. For predictive models with time to event

195  outcomes the number of participants experiencing the event should exceed 10 times the number

196  of degrees of freedom to ensure adequate sample size[18]. The number of dementia events in the

197  derivation cohort is 1 797 for men and 3 281 for women; therefore, the maximum number of

198  total degrees of freedom for each of the DemPoRT models is 179 and 328, respectively, which

199  we do not anticipate surpassing.

200

201  The validation cohorts will consist of approximately 4 600 males and 6 300 females, and 15 000

202  and 21 000 person-years of follow-up, respectively. Vergouwe et al[19] recommend a minimum of

203  100 events and 100 non-events for external validation studies. We expect approximately 225

204  events for men and 400 for women in our validation cohort.

205

206

207 **Analysis Plan**

208 The analysis plan was developed following guidelines by Harrell[18] and Steyerberg[20] after

209 accessing the derivation data set, but prior to model fitting or descriptive analyses involving

210 exposure-outcome associations. This was done to avoid Type 1 error introduced by data-driven

211 variable selection or model specification. Key considerations of our analysis approach include

212 full pre-specification of the predictor variables, use of flexible functions for continuous

213 predictors, and preserving statistical properties by avoiding data-driven variable selection

214 procedures. Analysis will be conducted using Harrell's Hmisc[21] package of functions in R[22] as

215 well as SAS v9.4.

216

217 This study protocol and the reporting of our model estimation results will be guided by the

218 TRIPOD statement for multivariable predictive models[23].

219

220 *Identification of Predictors*

221 Predictor variables were identified through review of existing predictive algorithms for

222 dementia[9,24–34] and comparison to available data collected in the CCHS. Variable inclusion was

223 informed by consultation with subject-matter experts and the project's advisory team, and

224 informed by our previous work developing predictive models for cardiovascular disease and life

225 expectancy[35,36].

226

227 Variables with more than 20% missing values, narrow distributions or insufficient variation were

228 excluded. Obvious cases of redundancy (e.g. alternate definitions of the same underlying

229 behaviour) were not included. A total of 32 predictor variables were identified: 7

230  sociodemographic, 3 general health, 9 behavioural, 7 functional, 5 health conditions and 1 design

231  variable (CCHS survey cycle). As the effect of dementia risk factors varies by sex, separate

232  models will be derived for men and women. Education, rather than individual income, was

233  selected as a predictor due to several concerns with income including lack of generalizability,

234  measurement error, stability over time and substantial missing values. Indicator variables for

235  smoking status were created to allow the inclusion of smoking pack-years as a continuous

236  predictor. The models will additionally include age interactions with the behavioural, functional

237  and health condition variables as the effect of these risk factors on dementia are expected to vary

238  with age. Detailed definitions and measurement of the predictor variables are presented in Table

239  1.

240

241  *Data Cleaning and Coding of Predictors*

242  Continuous variables will be inspected using boxplots and descriptive statistics to determine

243  values outside a plausible range. Values that are clearly erroneous will be corrected, where

244  possible, or set to missing. Continuous predictors with highly skewed distributions will be

245  truncated to the 99.5th percentile. Categorization of continuous variables will be avoided to

246  minimize the loss of predictive information. All data cleaning and coding will occur prior to

247  examining exposure-outcome associations.

248

249  *Missing Data*

250  As traditional complete cases analyses suffer from inefficiency, selection bias, and other

251  limitations[20], multiple imputation methods will be used to impute missing values using the

252  'aregImpute' function in the HMisc library[21]. This function simultaneously imputes missing

253    values using predictive mean matching and uses bootstrapping to take all aspects of uncertainty

254    in to account. The imputation model will consist of the full list of predictor variables, time to

255    event and censoring variables, as well as auxiliary variables that are not predictors, but may

256    nevertheless be useful in generating imputed values (e.g., income). The final model will be

257    estimated in each of five multiple imputation data sets and the results combined using the rules

258    developed by Rubin and Schenker[37] to account for imputation uncertainty.

259

260    *Model Specification*

261    Initial sex-specific main effects models will be fit using the pre-specified predictors and an initial

262    degree of freedom allocation for each predictor (Table 1). Decisions on initial degree of freedom

263    allocations were informed by the anticipated importance of each predictor and known dose-

264    response relationships with dementia. Continuous predictors will be flexibly modelled using

265    restricted cubic splines, with the knots placed at fixed quantiles of the distribution (e.g., $5^{th}$,

266    $27.5^{th}$, $50^{th}$, $72.5^{th}$, and $95^{th}$ centiles). Frequency distributions for categorical predictors will be

267    examined and categories with small numbers of respondents will be combined, with analysts

268    blinded to the number of events per category, to avoid instability in the regression analyses.

269    Ordinal variables will be specified as either linear terms or as categorical if the expected

270    association is more complex. Interactions will be restricted to linear terms. The initial model

271    specification, presented in Table 1, includes a total of 86 degrees of freedom (63 main, 23

272    interaction).

273

274    Partial association chi-square statistics for each predictor minus their degrees of freedom (to

275    level the playing field among predictors with varying degrees of freedom) will be plotted in

276 descending order. Variables with higher predictive potential will retain their initial degrees of

277 freedom, while predictors with lower predictive potential will be modeled as simple linear terms

278 or recoded by combining infrequent categories. This process of model specification does not

279 increase the Type I error rate because all predictors will be retained in the full model regardless

280 of their strength of association[18].

281

282 *Model Estimation*

283 The initial models will be estimated using competing risk Cox proportional hazards regression

284 with time to dementia ascertainment as the outcome and death as a competing risk. Alternative

285 model specifications, including flexible parametric models, will be considered after assessing the

286 validity of model assumptions. All predictors will be centered about their means. A formal check

287 of multicollinearity will be carried out using a variable clustering algorithm[18].

288

289 Proportional hazards models assume that the relative risk of the outcome between strata of

290 predictors and the baseline risk must be constant over time. Violation of this assumption has

291 been shown to produce biased results[38] although it has also been argued that the estimated

292 coefficients of time-varying variables can simply be interpreted as an average rather than

293 instantaneous hazard[39]. Plots of raw and smoothed scaled Schoenfeld residuals versus time for

294 each predictor will be assessed to test this assumption and identify non-proportionality. If a

295 violation of this assumption is identified we will consider addition of interaction terms between

296 the predictor and log-transformed time.

297

298 Although the risk of overfitting will be minimal due to pre-specification of the models and a

299 large sample size, the need for overfitting adjustment will be assessed. The degree of overfitting

300 will be estimated using the heuristic shrinkage estimator, based on the log likelihood ratio chi-

301 square statistic for the full model[40]. If shrinkage is <0.90, models will be adjusted for overfitting.

302

303 *Estimation of the Reduced Models*

304 Model pre-specification has advantages in limiting overfitting and spurious statistical

305 significance but can result in a final model that is overly complex, difficult to interpret, and

306 difficult to apply. Unnecessary predictor variables also distort the estimated effects of other

307 predictors making the model more computationally intensive. It is suggested that a more

308 parsimonious model that retains most of the prognostic information and performs as well as or

309 better than the full model can be derived without increasing the Type 1 error rate[18,41]. We will

310 identify a more parsimonious model using a stepdown procedure described by Ambler[41], which

311 involves deleting the variable that results in the smallest decrease in model $R^2$ until removal

312 leads to an $R^2$ that is less than a desired level. The reduced model will be evaluated against the

313 full model using Akaike's Information Criterion, and by examining the effect on discrimination

314 and calibration.

315

316 DemPoRT will be developed and validated using temporal split samples, however the final

317 regression coefficients will use the full data set to maximize follow-up duration. A cohort-

318 specific intercept and/or interaction term may be included in the final model if the derivation and

319 validation cohorts differ; otherwise, the final combined model will maintain the same predictors

320 and form as the derivation model.

321

*Assessment of Predictive Performance*

323    Predictive performance in the derivation and validation cohorts will be assessed and reported

324    using overall measures of predictive accuracy, discrimination and calibration. Accuracy will be

325    assessed with Nagelkerke's $R^2$ and the Brier score, and discrimination using the concordance

326    statistic. Model calibration is especially important in the development of prognostic models, as

327    probabilities of future risk are of primary interest[20,42,43]. Calibration will be assessed by

328    comparing the observed and predicted risk of dementia within vigintiles (20 groups of equal

329    frequency) of predicted risk with emphasis on visual inspection of plots rather than formal

330    statistical significance testing, which can be influenced by large sample sizes[19]. Calibration

331    slopes will be generated by regressing the outcome in the validation cohort on the predicted

332    dementia risk, reflecting the combined effect of overfitting to the derivation data as well as true

333    differences in effects of predictors. Deviation of the slope from 1 (perfect calibration) will be

334    tested using a Wald or likelihood ratio test. Calibration within predefined subgroups of

335    importance to clinicians and policy makers (e.g., age group, health behaviour, sociodemographic

336    groups and health conditions) will additionally be evaluated. The clinically relevant standard of

337    calibration was defined as less than 20% difference between observed and predicted estimates

338    within subgroups with a dementia prevalence of at least 5%. All model performance measures

339    will be calculated using the first of the multiply imputed data sets.

340

*Model Presentation*

342    The final regression model, derived from the combined sample of the derivation and validation

343    cohorts, will be presented using estimated hazard ratios and 95% confidence intervals, along

344  with results for the derivation and validation cohorts separately. We have found, however, this

345  usual presentation less meaningful when presenting complex models[35]. To allow interpretation of

346  the estimated effect of each predictor, model behavior will additionally be described using

347  interactive visual tools to display the shape of the effect of each predictor[44]. The regression

348  formula will also be published and used as the basis for web-based implementation.

349

350  **Analyses Beyond Initial Model Development**

351  We will conduct further analyses exploring the added predictive ability of novel risk factors that

352  were ascertained in single CCHS cycles (e.g. sedentary activity, cognitive stimulation, sleep

353  quality and duration), as well as risk factors that can be ascertained through linkage of additional

354  data sources and similar cohorts (e.g. detailed dietary consumption, lipid levels, blood pressure).

355  In addition, sensitivity analysis of the age at survey administration cutoff used for cohort creation

356  will be performed.

357

358  A second, causal model (DemPoRT-C) will also be created to assess the relative contribution of

359  lifestyle, socio-demographic and health factors to dementia incidence. Development will exclude

360  variables believed to be in the causal pathway of dementia occurrence (e.g., self-rated health and

361  functional measures) to reduce the attenuation of hazards from upstream risk factors, but will

362  otherwise be the same as in the predictive model. DemPoRT-C will be applied to the most recent

363  unlinked national CCHS survey.

364

365  **ETHICS AND MODEL DISSEMINATION**

366 The DemPoRT project advisory committee has been created to ensure that the models meet the

367 needs of knowledge users. This committee has worked with the study team to identify predictors

368 of dementia based on scientific and policy importance and will aid in the identification of

369 important target populations and the establishment of policy-relevant differences for calibration

370 studies.

371

372 DemPoRT results will be submitted for publication in peer-review journals and presented at

373 scientific meetings. A web-based individual-level calculator will be created if the models are

374 appropriate for individual use. Although DemPoRT emphasizes risk prediction at the population-

375 level, we have found that individual-level calculators are an effective engagement and translation

376 tool for both the general public and knowledge users.

377

378 **CONCLUSIONS**

379 To the best of our knowledge, DemPoRT will be the first population-based algorithm for

380 predicting and projecting dementia incidence. The DemPoRT models will produce estimates of

381 future dementia burden that we believe will be more accurate than existing estimates, will assess

382 the contribution of specific risk factors to the population risk, and identify groups at high risk of

383 developing dementia. Although a rigorous approach to model development will be used, further

384 validation will be needed to assess generalizability, and calibration will be required for

385 application in other jurisdictions.

386

387 **CONTRIBUTIONS**

388 SF drafted and revised the manuscript, and contributed to study design and protocol

389 development. NM contributed to study design, protocol development and provided

390 data/statistical support. AH, MT, DM and GH contributed to the design of the study and protocol

391 development. PT is the lead investigator of the study and was responsible for the conception of

392 the project, the grant application, study design and protocol development. All authors provided

393 critical reviews of the manuscript and approved the final version.

394

407

408 **COMPETING INTERESTS**

409 None declared.

410

411 **ETHICS APPROVAL**

412 Research ethics approval has been obtained from the Ottawa Health Science Network Research

413 Ethics Board.

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

**REFERNCES**

1.  Brookmeyer, R., Gray, S. & Kawas, C. Projections of Alzheimer's disease in the United States and the public health impact of delaying disease onset. *Am. J. Public Health* **88,** 1337–1342 (1998).

2.  Brayne, C. The elephant in the room - healthy brains in later life, epidemiology and public health. *Nat. Rev. Neurosci.* **8,** 233–9 (2007).

3.  Norton, S., Matthews, F. E. & Brayne, C. A commentary on studies presenting projections of the future prevalence of dementia. *BMC Public Health* **13,** 1 (2013).

4.  Sloane, P. D. *et al.* The public health impact of Alzheimer's disease, 2000-2050: potential implication of treatment advances. *Annu. Rev. Public Health* **23,** 213–31 (2002).

5.  Mura, T., Dartigues, J. F. & Berr, C. How many dementia cases in France and Europe? Alternative projections and scenarios 2010-2050. *Eur. J. Neurol.* **17,** 252–259 (2010).

6.  Hebert, L. E., Scherr, P. A., Bienias, J. L., Bennett, D. A. & Evans, D. A. Alzheimer disease in the US population: prevalence estimates using the 2000 census. *Arch. Neurol.* **60,** 1119–22 (2003).

7.  Brookmeyer, R., Johnson, E., Ziegler-Graham, K. & Arrighi, H. M. Forecasting the global burden of Alzheimer's disease. *Alzheimer's Dement.* **3,** 186–191 (2007).

8.  Dewey, M. E. & Chen, C.-M. Neurosis and mortality in persons aged 65 and over living in the community: a systematic review of the literature. *Int. J. Geriatr. Psychiatry* **19,** 554–7 (2004).

9.  Norton, S., Matthews, F. E., Barnes, D. E., Yaffe, K. & Brayne, C. Potential for primary prevention of Alzheimer's disease: An analysis of population-based data. *Lancet Neurol.* **13,** 788–794 (2014).

457    10.    Joly, P. *et al.* Prevalence projections of chronic diseases and impact of public health

458            intervention. *Biometrics* **69,** 109–17 (2013).

459    11.    Lee, Y. The recent decline in prevalence of dementia in developed countries: implications

460            for prevention in the Republic of Korea. *J. Korean Med. Sci.* **29,** 913–8 (2014).

461    12.    Barnes, D. E. & Yaffe, K. The projected effect of risk factor reduction on Alzheimer's

462            disease prevalence. *Lancet Neurol.* **10,** 819–828 (2011).

463    13.    Manuel, D. G. *et al.* Projections of preventable risks for cardiovascular disease in Canada

464            to 2021: a microsimulation modelling approach. *C. open* **2,** E94–E101 (2014).

465    14.    Béland, Y. Canadian community health survey--methodological overview. *Heal. reports /*

466            *Stat. Canada, Can. Cent. Heal. Inf. = Rapp. sur la sant?? / Stat. Canada, Cent. Can.*

467            *d'information sur la sant??* **13,** 9–14 (2002).

468    15.    Jaakkimainen, R. L. *et al.* Identification of Physician-Diagnosed Alzheimer's Disease and

469            Related Dementias in Population-Based Administrative Data: A Validation Study Using

470            Family Physicians' Electronic Medical Records. *J. Alzheimer's Dis.* **54,** 337–349 (2016).

471    16.    Connolly, A., Gaehl, E., Martin, H., Morris, J. & Purandare, N. Underdiagnosis of

472            dementia in primary care: Variations in the observed prevalence and comparisons to the

473            expected prevalence. *Aging Ment. Health* **15,** 978–984 (2011).

474    17.    Kosteniuk, J. G. *et al.* Incidence and prevalence of dementia in linked administrative

475            health data in Saskatchewan, Canada: a retrospective cohort study. *BMC Geriatr.* **15,** 73

476            (2015).

477    18.    Harrell, F. E. *Regression Modeling Strategies with applicaitons to linear models, logistic*

478            *regression and survival analysis.* (Springer, 2001).

479    19.    Vergouwe, Y., Steyerberg, E. W., Eijkemans, M. J. C. & Habbema, J. D. F. Substantial

480    effective sample sizes were required for external validation studies of predictive logistic

481    regression models. *J. Clin. Epidemiol.* **58,** 475–483 (2005).

482    20.    Steyerberg, E. W. *Clinical Prediction Models*. (Springer, 2009).

483    21.    Harrell, F. E. *Hmisc: Harrell Miscellaneous. R package version 4.0-2*. (2016).

484    22.    R Core Team. R: A language and environment for statistical computing. (2016).

485    23.    Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. M. Transparent reporting of

486    a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The

487    TRIPOD Statement. *Eur. Urol.* **67,** 1142–1151 (2015).

488    24.    Walters, K. *et al.* Predicting dementia risk in primary care: development and validation of

489    the Dementia Risk Score using routinely collected data. *BMC Med.* **14,** 1–12 (2016).

490    25.    Anstey, K. J. *et al.* A self-report risk index to predict occurrence of dementia in three

491    independent cohorts of older adults: The ANU-ADRI. *PLoS One* **9,** (2014).

492    26.    Exalto, L. G. *et al.* Midlife risk score for the prediction of dementia four decades later.

493    *Alzheimer's Dement.* **10,** 562–570 (2014).

494    27.    Jessen, F. *et al.* Prediction of dementia in primary care patients. *PLoS One* **6,** e16852

495    (2011).

496    28.    Reitz, C. *et al.* A summary risk score for the prediction of Alzheimer disease in elderly

497    persons. *Arch. Neurol.* **67,** 835–41 (2010).

498    29.    Barnes, D. E. *et al.* Predicting risk of dementia in older adults: The late-life dementia risk

499    index. *Neurology* **73,** 173–179 (2009).

500    30.    Kivipelto, M. *et al.* Risk score for the prediction of dementia risk in 20 years among

501    middle aged people: a longitudinal, population-based study. *Lancet Neurol.* **5,** 735–741

502    (2006).

503    31.    Song, X., Mitnitski, A. & Rockwood, K. Nontraditional risk factors combine to predict

504           Alzheimer disease and dementia. *Neurology* **77,** 227–234 (2011).

505    32.    Tierney, M. C., Moineddin, R. & McDowell, I. Prediction of all-cause dementia using

506           neuropsychological tests within 10 and 5 years of diagnosis in a community-based sample.

507           *J. Alzheimer's Dis.* **22,** 1231–1240 (2010).

508    33.    Exalto, L. G. *et al.* Risk score for prediction of 10 year dementia risk in individuals with

509           type 2 diabetes: A cohort study. *Lancet Diabetes Endocrinol.* **1,** 183–190 (2013).

510    34.    Chary, E. *et al.* Short-versus long-term prediction of dementia among subjects with low

511           and high educational levels. *Alzheimer's Dement.* **9,** 562–571 (2013).

512    35.    Taljaard, M. *et al.* Cardiovascular Disease Population Risk Tool (CVDPoRT): predictive

513           algorithm for assessing CVD risk in the community setting. A study protocol. *BMJ Open*

514           **4,** e006701 (2014).

515    36.    Manuel, D. G. *et al.* Measuring Burden of Unhealthy Behaviours Using a Multivariable

516           Predictive Approach: Life Expectancy Lost in Canada Attributable to Smoking, Alcohol,

517           Physical Inactivity, and Diet. *PLoS Med.* **13,** 1–27 (2016).

518    37.    Rubin, D. B. & Schenker, N. Multiple imputation in health-care databases: an overview

519           and some applications. *Stat. Med.* **10,** 585–98 (1991).

520    38.    Therneau, T. M., Grambsch, P. M. & Fleming, T. R. Martingale-based residuals for

521           survival models. *Biometrika* **77,** 147–160 (1990).

522    39.    Allison, P. D. *Survival analysis using SAS : A practical guide*. (SAS Institute, 2010).

523    40.    Van Houwelingen, J. C. & Le, C. S. Predictive value of statistical models. *Stat Med* **9,**

524           1303–1325 (1990).

525    41.    Ambler, G., Brady, A. R. & Royston, P. Simplifying a prognostic model: A simulation

526 study based on clinical data. *Stat. Med.* **21,** 3803–3822 (2002).

527 42. Cook, N. R. Statistical evaluation of prognostic versus diagnostic models: beyond the

528 ROC curve. *Clin. Chem.* **54,** 17–23 (2008).

529 43. Cook, N. R. Comment: Measures to summarize and compare the predictive capacity of

530 markers. *Int. J. Biostat.* **6,** Article 22; discussion Article 25 (2010).

531 44. Krause, J., Perer, A. & Bertini, E. Using Visual Analytics to Interpret Predictive Machine

532 Learning Models. *arXiv Prepr. arXiv1606.05685.* (2016).

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550    Table 1. Pre-specification of predictor variables for DemPoRT with initial degrees of freedom (df) allocation

| Variable | Scale | Initial Variable Specification | df |
|---|---|---|---|
| **Socio-demographic Factors** | | | |
| Age | Continuous | **5 knot spline:** Valid range: 55-102 (male), 55-101 (female) | 4 |
| Sex | Categorical | **Stratified:** Male; Female | NA |
| Ethnicity | Categorical | **7 categories:** Caucasian; African-American; Chinese; Aboriginal; Japanese/Korean/South East Asian/Filipino; Other/Multiple origin/Unknown/Latin American; South Asian/Arab/West Asian | 6 |
| Immigrant | Dichotomous | Yes; No | 1 |
| Education | Categorical | **4 categories:** Less than secondary school; Secondary school graduation; Some postsecondary; Postsecondary graduation | 3 |
| Marital Status | Categorical | **4 categories:** Now married/Common-law; Separated/Divorced; Widowed; Single | 3 |
| Neighborhood Social and Material Deprivation (Pampalon et al. 2009) | Ordinal | **3 categories:** Low (1st or 2nd quintile); High 4th or 5th quintile; Moderate (3rd quintile) | 2 |
| **General Health** | | | |
| Sense of belonging to local community | Ordinal | **4 categories:** Very strong; Somewhat strong; Somewhat weak; Very weak | 3 |
| Self-perceived stress | Ordinal | **5 categories:** Not at all stressful; Not very stressful; A bit stressful; Quite a bit stressful; Extremely stressful | 4 |
| Self-rated health | Ordinal | **5 categories:** Poor; Fair; Good; Very Good; Excellent | 4 |
| **Health Behaviors** | | | |
| Pack years of smoking | Continuous | **3 knot spline:** Valid range: 0-112 (male), 0-78 (female) | 2 |
| Smoking status | Categorical | **4 categories:** Non-smoker; Current smoker; Former smoker quit <5 years ago; Former smoker quit >5 years ago | 3 |
| Alcohol consumption (number of drinks last week) | Continuous | **3 knot spline:** Valid range: 0-50 (male), 0-24 (female) | 2 |
| Former drinker | Dichotomous | Yes; No | 1 |
| Consumption of fruit, salad, carrot and other vegetables (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-48 (male), 0-31 (female) | 2 |
| Potato consumption (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-2 | 2 |
| Juice consumption (average daily consumption | Continuous | **3 knot spline:** Valid range: 0-6 (male), 0-5 (female) | 2 |
| Leisure physical activity (average daily METs (kcal/kg/day)) | Continuous | **3 knot spline:** Valid range: 0-16 (male), 0-12 (female) | 2 |
| **Functional Measures** | | | |
| Personal hygiene and care | Dichotomous | Does not need help; Needs help | 1 |
| Locomotion in the home | Dichotomous | Does not need help; Needs help | 1 |
| Meal preparation | Dichotomous | Does not need help; Needs help | 1 |
| Running errands | Dichotomous | Does not need help; Needs help | 1 |
| Ordinary housework | Dichotomous | Does not need help; Needs help | 1 |
| Heavy housework | Dichotomous | Does not need help; Needs help | 1 |
| Finances | Dichotomous | Does not need help; Needs help | 1 |
| **Health Conditions** | | | |
| Heart disease | Dichotomous | Yes; No | 1 |
| Stroke | Dichotomous | Yes; No | 1 |
| Diabetes | Dichotomous | Yes; No | 1 |
| Mood disorder | Dichotomous | Yes; No | 1 |
| High blood pressure | Dichotomous | Yes; No | 1 |
| Body mass index | Continuous | **3 knot spline:** Valid range: 10-44 (male), 10-47 (female) | 2 |
| **Design** | | | |
| Survey year | Ordinal | **4 categories:** 2000/01, 2002/03, 2004/05, 2006/07 | 3 |

551    DemPoRT, Dementia Population Risk Tool; df, degrees of freedom; MET, metabolic equivalent task

# TRIPOD Checklist: Prediction Model Development and Validation

| Section/Topic | Item | | Checklist Item | Page |
|---|---|---|---|---|
| **Title and abstract** | | | | |
| Title | 1 | D;V | Identify the study as developing and/or validating a multivariable prediction model, the target population, and the outcome to be predicted. | 1 |
| Abstract | 2 | D;V | Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions. | 3 |
| **Introduction** | | | | |
| Background and objectives | 3a | D;V | Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models. | 5,6 |
| | 3b | D;V | Specify the objectives, including whether the study describes the development or validation of the model or both. | 7 |
| **Methods** | | | | |
| Source of data | 4a | D;V | Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable. | 7,8 |
| | 4b | D;V | Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up. | 7,8 |
| Participants | 5a | D;V | Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres. | 7,8 |
| | 5b | D;V | Describe eligibility criteria for participants. | 7-9 |
| | 5c | D;V | Give details of treatments received, if relevant. | NA |
| Outcome | 6a | D;V | Clearly define the outcome that is predicted by the prediction model, including how and when assessed. | 8,9 |
| | 6b | D;V | Report any actions to blind assessment of the outcome to be predicted. | 10 |
| Predictors | 7a | D;V | Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured. | 10,11, 25 |
| | 7b | D;V | Report any actions to blind assessment of predictors for the outcome and other predictors. | 10 |
| Sample size | 8 | D;V | Explain how the study size was arrived at. | 7-9 |
| Missing data | 9 | D;V | Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method. | 11,12 |
| Statistical analysis methods | 10a | D | Describe how predictors were handled in the analyses. | 12,13, 25 |
| | 10b | D | Specify type of model, all model-building procedures (including any predictor selection), and method for internal validation. | 13,14 |
| | 10c | V | For validation, describe how the predictions were calculated. | 15 |
| | 10d | D;V | Specify all measures used to assess model performance and, if relevant, to compare multiple models. | 15 |
| | 10e | V | Describe any model updating (e.g., recalibration) arising from the validation, if done. | NA |
| Risk groups | 11 | D;V | Provide details on how risk groups were created, if done. | 15 |
| Development vs. validation | 12 | V | For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors. | 15,16 |
| **Results** | | | | |
| Participants | 13a | D;V | Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful. | 9 |
| | 13b | D;V | Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome. | NA |
| | 13c | V | For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome). | NA |
| Model development | 14a | D | Specify the number of participants and outcome events in each analysis. | NA |
| | 14b | D | If done, report the unadjusted association between each candidate predictor and outcome. | NA |
| Model specification | 15a | D | Present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point). | NA |
| | 15b | D | Explain how to the use the prediction model. | NA |
| Model performance | 16 | D;V | Report performance measures (with CIs) for the prediction model. | NA |
| Model-updating | 17 | V | If done, report the results from any model updating (i.e., model specification, model performance). | NA |
| **Discussion** | | | | |
| Limitations | 18 | D;V | Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data). | 4,17 |
| Interpretation | 19a | V | For validation, discuss the results with reference to performance in the development data, and any other validation data. | NA |
| | 19b | D;V | Give an overall interpretation of the results, considering objectives, limitations, results from similar studies, and other relevant evidence. | NA |
| Implications | 20 | D;V | Discuss the potential clinical use of the model and implications for future research. | 5-7 |
| **Other information** | | | | |
| Supplementary information | 21 | D;V | Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets. | 4,17 |
| Funding | 22 | D;V | Give the source of funding and the role of the funders for the present study. | 18 |

*Items relevant only to the development of a prediction model are denoted by D, items relating solely to a validation of a prediction model are denoted by V, and items relating to both are denoted D;V. We recommend using the TRIPOD Checklist in conjunction with the TRIPOD Explanation and Elaboration document.
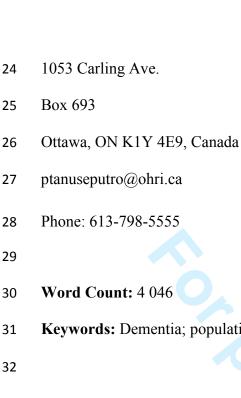
# BMJ Open

## Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm Assessing Dementia Risk in the Community

SCHOLARONE™
Manuscripts

1  **Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm**

2  **Assessing Dementia Risk in the Community**

3

4  Stacey Fisher, MSc[1,2,3] stacey.fisher@uottawa.ca

5  Amy Hsu, PhD[1,2] ahsu@toh.ca

6  Nassim Mojaverian, MSc[2] namojaverian@ohri.ca

7  Monica Taljaard, PhD[1,3] mtaljaard@ohri.ca

8  Gregory Huyer, PhD[4] ghuye047@uottawa.ca

9  Doug Manuel, MD[1,2,3,5,6] dmanuel@ohri.ca

10  Peter Tanuseputro, MD[1,2,6,7,8] ptanuseputro@ohri.ca

11

12  [1] Ottawa Hospital Research Institute, Ottawa, Ontario, Canada

13  [2] Institute for Clinical Evaluative Sciences, Ottawa, Ontario, Canada

14  [3] Department of Epidemiology and Community Medicine, University of Ottawa, Ottawa,

15  Ontario, Canada

16  [4] Telfer School of Management, University of Ottawa, Ottawa, Ontario, Canada

17  [5] Statistics Canada, Ottawa, Ontario, Canada

18  [6] Department of Family Medicine, University of Ottawa, Ottawa, Ontario, Canada

19  [7] Department of Medicine, University of Ottawa, Ottawa, Ontario, Canada

20  [8] Bruyère Research Institute, Ottawa, Ontario, Canada

21

22  Corresponding author:

23  Dr. Peter Tanuseputro

24 1053 Carling Ave.

25 Box 693

26 Ottawa, ON K1Y 4E9, Canada

27 ptanuseputro@ohri.ca

28 Phone: 613-798-5555

29

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47    **ABSTRACT**

48    **Introduction:** The burden of disease from dementia is a growing global concern as incidence

49    increases dramatically with age and average life expectancy has been increasing around the

50    world. Planning for an aging population requires reliable projections of dementia prevalence;

51    however, existing population projections are simple and have poor predictive accuracy. The

52    Dementia Population Risk Tool (DemPoRT) will predict incidence of dementia in the population

53    setting using multivariable modeling techniques, and will be used to project dementia

54    prevalence.

55    **Methods and Analysis:** The derivation cohort will consist of elderly Ontario respondents of the

56    Canadian Community Health Survey (CCHS) (2001, 2003, 2005, 2007; 18 764 males and 25 288

57    females). Pre-specified predictors include sociodemographic, general health, behavioral,

58    functional and health condition variables. Incident dementia will be identified through individual

59    linkage of survey respondents to population-level administrative health care databases (1 797 and

60    3 281 events, and 117 795 and 166 573 person-years of follow-up, for males and females,

61    respectively until March 31, 2014). Using time of first dementia capture as the primary outcome

62    and death as a competing risk, sex-specific proportional hazards regression models will be

63    estimated. The 2008/2009 CCHS survey will be used for validation (approximately 4 600 males

64    and 6 300 females). Overall calibration and discrimination will be assessed as well as calibration

65    within predefined subgroups of importance to clinicians and policy makers.

66    **Ethics and Dissemination:** Research ethics approval has been granted by the Ottawa Health

67    Science Network Research Ethics Board. DemPoRT results will be submitted for publication in

68    peer-review journals and presented at scientific meetings. The algorithm will be assessable

69    online for both population and individual uses.

70      **Trial Registration Number:** ClinicalTrials.gov NCT03155815.

71

72      **STRENGTHS AND LIMITATIONS**

73      - The Dementia Population Risk Tool (DemPoRT) will be developed and validated using

74          predictors from large population-based community health surveys that are individually

75          linked to routinely-collected health administration data in Ontario. To our knowledge,

76          DemPoRT will be the first algorithm designed to predict and project dementia incidence

77          at the population-level.

78      - DemPoRT will be used to produce improved estimates of future dementia burden, will

79          assess the contribution of specific risk factors to the population risk, and will identify

80          population subgroups at high risk of developing dementia. This information will be used

81          by policymakers to prepare for and reduce dementia impact.

82      - The analysis plan and predictors have been fully pre-specified to limit the risk of over-

83          fitting and improve the quality of predictions.

84      - Detailed cognitive testing to ascertain dementia diagnoses is preferable over the use of

85          administrative data, however this is not available or feasible at the population level.

86      - Although a rigorous approach to model development will be used, further validation will

87          be needed to assess generalizability, and calibration will be required for application in

88          other jurisdictions.

89

90

91

92

93 **INTRODUCTION**

94 The burden of disease from dementia is a growing global concern as incidence increases

95 dramatically with age and average life expectancy has been increasing around the world[1,2].

96 Planning for an aging population requires reliable projections of dementia burden and the

97 implications for resource requirements. Existing population-level projections for dementia,

98 however, are overly simplistic and likely inaccurate[3].

99

100 **Limitations of Current Dementia Projection Methodology**

101 Almost all existing dementia projections have used extrapolation and macrosimulation methods,

102 which are simplistic and make assumptions that may not hold true into the future[3]. Most

103 extrapolations simply apply current age- and sex-specific prevalence estimates of dementia to

104 future population projections. Macrosimulations typically use estimates of dementia incidence

105 and mortality, stratified by age and sex, to simulate disease prevalence as the population ages[1,4–

106 6]. Projections from extrapolations incorrectly assume that the risk of mortality among those with

107 and without dementia are equivalent[7,8], and both methods assume that the age and sex-specific

108 prevalence of dementia risk factors will not change with time. The assumption of stable risk

109 factor prevalence is widely thought to be the major source of error in existing dementia

110 projections[3,9–11].

111

112 Changing trends of dementia risk factors has the potential to have a dramatic impact on dementia

113 prevalence estimates, as up to 50% of dementia cases have been attributed to modifiable

114 factors[9,12], and the prevalence of several factors has been projected to change significantly in the

115 near future. For example, the population prevalence of diabetes and obesity in Canada has been

116  projected to increase, while smoking, hypertension and dyslipidemia have been projected to

117  decline[13]. Consideration of risk factor prevalence is therefore important to improve the accuracy

118  of dementia projections, and simple extrapolations and macrosimulations are often inadequate.

119

120  **Predictive Multivariable Modeling of Dementia Incidence**

121  Population-based **predictive risk algorithms** examine the effect of risk factors on dementia

122  incidence, and can be used for dementia burden projection. Population-based data that contain

123  detailed risk factor information, such as health surveys, are linked at the individual-level to

124  administrative data that capture dementia development. A multivariable model of dementia

125  incidence is derived, validated against external data, and predictive performance is assessed.

126  Once developed, the algorithm can be used to project disease incidence and prevalence. To

127  obtain prevalence projections, the algorithm can be integrated in to a microsimulation model

128  such as Statistics Canada's Population Health Models (POHEM). POHEM dynamically models

129  individual life trajectories of a population representative of Canada including births, deaths and

130  migration, disease incidence and progression, and exposure to risk factors, facilitating detailed

131  examination of the influence of changing risk factor prevalence on future dementia prevalence.

132

133  Predictive risk algorithms can also be used to describe the risk of dementia in the population,

134  assess the contribution of specific risk factors to the population risk, identify high-risk groups,

135  and evaluate risk reduction strategies.

136

137  **Existing Dementia Prediction Models**

138  Many models have been developed to predict risk of dementia[14–26], most with the primary goal of

139  identifying individuals in the clinical setting at high risk. They have varying discriminative

140  ability (c-statistics ranging from 0.49[16] to 0.89[17]) and have generally been derived from small

141  samples, rarely including more than a few thousand individuals. Existing models are therefore

142  simplistic, including few predictors and rarely including interaction or non-linear terms Existing

143  models thus facilitates understanding and use by physicians in clinical practice, but limits

144  discriminatory ability and predictive accuracy. Walters et al[26] developed an algorithm for

145  predicting 5-year dementia risk among individuals 60-79 years of age in the United Kingdom

146  using an enormous derivation dataset of 800 000 individuals, and a simple model. The derivation

147  model had a c-statistic of 0.84 (95% CI: 0.81, 0.87), but a low positive predictive value at most

148  risk thresholds, and therefore is poor at identifying those at high risk of dementia. Additionally,

149  as most dementia risk models are intended for use in the clinical setting, many include results

150  from neuropsycological tests[17–23], MRI findings[18] and APOE genotype[18,24,25]. The inclusion of

151  these variables, however, limits the application of these models as these variables are not

152  available at the population-level.

153

154  The objective of this study is to develop and validate the Dementia Population Risk Tool

155  (DemPoRT) algorithm to predict dementia incidence in the population setting. This will be done

156  using multivariable modelling techniques, linking self-reported risk factors captured by a

157  population-based health survey in Canada with administrative databases across healthcare sectors

158  that capture healthcare diagnosed dementia. DemPoRT will be developed with a using a large

159  population-based dataset using only variables that are available at the population-level, allowing

160  for population-level application. DemPoRT will also utilize many methodological improvements

161     over existing models. This protocol pre-specifies the predictor variables and analytic plan for

162     model development, reducing the potential for overfitting and bias, and improving transparency.

163     Interaction terms and flexible functions for continuous predictors will be investigated, increasing

164     potential discriminative ability. The pre-specified analytic plan avoids data-driven variable

165     selection procedures, further reducing the potential for bias.

166

167     To our knowledge, the DemPoRT predictive model will be the first algorithm designed to predict

168     and project dementia incidence at the population-level. It will be used to estimate the future

169     burden of dementia using techniques that consider changes in risk factor prevalence and will

170     identify modifiable risk factors that can be targeted by individuals, clinicians and policy makers

171     to reduce the burden of dementia.

172

173     **METHODS AND ANALYSIS**

174     **Study Design**

175     Two DemPoRT models, one for males and females, will be derived and validated using

176     population-based data in Ontario, Canada, a multicultural province with 13.6 million residents.

177     Predictors will be obtained from the Canadian Community Health Surveys (CCHS), and

178     outcomes (i.e., diagnosis of dementia) will be obtained from routinely-collected health care data.

179

180     The derivation cohorts will consist of eligible respondents of the 2001, 2003, 2005 and 2007

181     CCHS (Cycles 1.1, 2.1, 3.1 and 4.1), while validation cohorts will consist of respondents to the

182     2008/2009 cycle. The CCHS is a national, cross-sectional survey developed by Statistics Canada

183     to collect information related to health and health care utilization of the Canadian population.

184    The survey has a multistage stratified cluster design that represents approximately 98% of the

185    Canadian population aged 12 years and over and attained an average response rate of 79% over

186    the study period. The CCHS is conducted through telephone and in-person interviews, and all

187    responses are self-reported. The details of survey methodology have been published elsewhere[27].

188    Survey respondents will be excluded if they are less than 55 years of age at survey

189    administration, self-reported a history of dementia, or are not eligible for Ontario's universal

190    health insurance. If a respondent was included in more than one CCHS cycle, only their earliest

191    survey response will be used.

192

193    **Outcome**

194    Survey respondents diagnosed with dementia will be identified through individual linkage to

195    several population-based administrative databases at the Institute for Clinical Evaluative

196    Sciences (ICES). Dementia case ascertainment is based on a validated definition: 1 hospital

197    record OR 3 physician claim records at least 30 days apart within a 2-year period OR a

198    dispensing record for a cholinesterase inhibitor from Ontario Drug Benefit (ODB). This

199    definition has a sensitivity of 79.3% and a specificity of 99.1% when validated against

200    emergency medical record (EMR) data[28]. Due to known underdiagnosis of dementia[29,30], we will

201    supplement this definition by adding survey respondents with dementia codes captured on home

202    care and long-term care assessments (dementia flag AND Cognitive Performance Scale [CPS]

203    score $\geq 2$) using the Resident Assessment Instrument-Home Care (RAI-HC) database and the

204    Continuing Care Reporting System (CCRS), respectively. We have found this addition adds

205    substantially (approximately 18%) to the number of dementia cases captured.

206

207 Survey respondents with dementia will be excluded if they meet the criteria for dementia within

208 2 years of survey administration (to remove potentially prevalent cases) or are younger than 65

209 years of age at the time of dementia diagnosis (to exclude early onset dementia which likely has

210 a different set of risk factors). Eligible survey respondents will be followed from the date of

211 survey administration or age 65, whichever came later, until the earliest date of: dementia

212 ascertainment, death (defined as competing risk), loss to follow-up (defined as loss of healthcare

213 eligibility) or end of study (March 31, 2014).

214

215 **Sample Size**

216 The male and female derivation cohorts consist of 18 764 and 25 288 respondents, and 117 795

217 and 166 573 person-years of follow-up, respectively. For predictive models with time to event

218 outcomes the number of participants experiencing the event should exceed 10 times the number

219 of degrees of freedom to ensure adequate sample size[31]. The number of dementia events in the

220 derivation cohort is 1 797 for men and 3 281 for women; therefore, the maximum number of

221 total degrees of freedom for each of the DemPoRT models is 179 and 328, respectively, which

222 we do not anticipate surpassing.

223

224 The validation cohorts will consist of approximately 4 600 males and 6 300 females, and 15 000

225 and 21 000 person-years of follow-up, respectively. Vergouwe *et al*[32] recommend a minimum of

226 100 events and 100 non-events for external validation studies. We expect approximately 225

227 events for men and 400 for women in our validation cohort.

228

229 **Analysis Plan**

230   The analysis plan was developed following guidelines by Harrell[31] and Steyerberg[33] after

231   accessing the derivation data set, but prior to model fitting or descriptive analyses involving

232   exposure-outcome associations. This was done to avoid Type 1 error introduced by data-driven

233   variable selection or model specification. Key considerations of our analysis approach include

234   full pre-specification of the predictor variables, use of flexible functions for continuous

235   predictors, and preserving statistical properties by avoiding data-driven variable selection

236   procedures. Analysis will be conducted using Harrell's Hmisc[34] package of functions in R[35] as

237   well as SAS v9.4.

238

239   This study protocol and the reporting of our model estimation results will be guided by the

240   TRIPOD statement for multivariable predictive models[36].

241

242   *Identification of Predictors*

243   Predictor variables were identified through review of existing predictive algorithms for

244   dementia[9,16,18–22,24–26,37,38] and comparison to available data collected in the CCHS. Variable

245   inclusion was informed by consultation with subject-matter experts and the project's advisory

246   team, and informed by our previous work developing predictive models for cardiovascular

247   disease and life expectancy[39,40].

248

249   Variables with narrow distributions or insufficient variation were excluded. Obvious cases of

250   redundancy (e.g. alternate definitions of the same underlying behavior) were not included. A

251   total of 32 predictor variables were identified: 7 sociodemographic, 3 general health, 9

252   behavioral, 7 functional, 5 health conditions and 1 design variable (CCHS survey cycle). As the

253  effect of dementia risk factors varies by sex, separate models will be derived for men and

254  women. Education, rather than individual income, was selected as a predictor due to several

255  concerns with income including lack of generalizability, measurement error, stability over time

256  and substantial missing values. Neighborhood social and material deprivation is captured using

257  Pampalon's deprivation index[41]. Indicator variables for smoking status were created to allow the

258  inclusion of smoking pack-years as a continuous predictor. The models will additionally include

259  age interactions with the behavioral, functional and health condition variables as the effect of

260  these risk factors on dementia are expected to vary with age. Detailed definitions and

261  measurement of the predictor variables are presented in Table 1.

262

263  *Data Cleaning and Coding of Predictors*

264  Continuous variables will be inspected using boxplots and descriptive statistics to determine

265  values outside a plausible range. Values that are clearly erroneous will be corrected, where

266  possible, or set to missing. Continuous predictors with highly skewed distributions will be

267  truncated to the 99.5th percentile. Categorization of continuous variables will be avoided to

268  minimize the loss of predictive information. All data cleaning and coding will occur prior to

269  examining exposure-outcome associations.

270

271  *Missing Data*

272  As traditional complete cases analyses suffer from inefficiency, selection bias, and other

273  limitations[33], multiple imputation methods will be used to impute missing values using the

274  'aregImpute' function in the HMisc library[34]. This function simultaneously imputes missing

275  values using predictive mean matching and uses bootstrapping to take all aspects of uncertainty

276  in to account. The imputation model will consist of the full list of predictor variables, time to

277  event and censoring variables, as well as auxiliary variables that are not predictors, but may

278  nevertheless be useful in generating imputed values (e.g., income). The final model will be

279  estimated in each of five multiple imputation data sets and the results combined using the rules

280  developed by Rubin and Schenker[42] to account for imputation uncertainty.

281

*Model Specification*

283  Initial sex-specific main effects models will be fit using the pre-specified predictors and an initial

284  degree of freedom allocation for each predictor (Table 1). Decisions on initial degree of freedom

285  allocations were informed by the anticipated importance of each predictor and known dose-

286  response relationships with dementia. Continuous predictors will be flexibly modelled using

287  restricted cubic splines, with the knots placed at fixed quantiles of the distribution (e.g., $5^{th}$,

288  $27.5^{th}$, $50^{th}$, $72.5^{th}$, and $95^{th}$ centiles). Frequency distributions for categorical predictors will be

289  examined and categories with small numbers of respondents will be combined, with analysts

290  blinded to the number of events per category, to avoid instability in the regression analyses.

291  Ordinal variables will be specified as either linear terms or as categorical if the expected

292  association is more complex. Interactions will be restricted to linear terms. The initial model

293  specification, presented in Table 1, includes a total of 86 degrees of freedom (63 main, 23

294  interaction).

295

296  Partial association chi-square statistics for each predictor minus their degrees of freedom (to

297  level the playing field among predictors with varying degrees of freedom) will be plotted in

298  descending order. Variables with higher predictive potential will retain their initial degrees of

299   freedom, while predictors with lower predictive potential will be modeled as simple linear terms

300   or recoded by combining infrequent categories. This process of model specification does not

301   increase the Type I error rate because all predictors will be retained in the full model regardless

302   of their strength of association[31].

303

304   *Model Estimation*

305   The initial models will be estimated using competing risk Cox proportional hazards regression

306   with time to dementia ascertainment as the outcome and death as a competing risk. Alternative

307   model specifications, including subdistribution hazard and flexible parametric models, will be

308   considered. All predictors will be centered about their means. A formal check of

309   multicollinearity will be carried out using a variable clustering algorithm[31].

310

311   Proportional hazards models assume that the relative risk of the outcome between strata of

312   predictors and the baseline risk must be constant over time. Violation of this assumption has

313   been shown to produce biased results[43] although it has also been argued that the estimated

314   coefficients of time-varying variables can simply be interpreted as an average rather than

315   instantaneous hazard[44]. Plots of raw and smoothed scaled Schoenfeld residuals versus time for

316   each predictor will be assessed to test this assumption and identify non-proportionality. If a

317   violation of this assumption is identified we will consider addition of interaction terms between

318   the predictor and log-transformed time.

319

320   Although the risk of overfitting will be minimal due to pre-specification of the models and a

321   large sample size, the need for overfitting adjustment will be assessed. The degree of overfitting

14

322 will be estimated using the heuristic shrinkage estimator, based on the log likelihood ratio chi-

323 square statistic for the full model[45]. If shrinkage is <0.90, models will be adjusted for overfitting.

324

*Estimation of the Reduced Models*

326 Model pre-specification has advantages in limiting overfitting and spurious statistical

327 significance but can result in a final model that is overly complex, difficult to interpret, and

328 difficult to apply. Unnecessary predictor variables also distort the estimated effects of other

329 predictors making the model more computationally intensive. It is suggested that a more

330 parsimonious model that retains most of the prognostic information and performs as well as or

331 better than the full model can be derived without increasing the Type 1 error rate[31,46]. We will

332 identify a more parsimonious model using a stepdown procedure described by Ambler[46], which

333 involves deleting the variable that results in the smallest decrease in model $R^2$ until removal

334 leads to an $R^2$ that is less than a desired level. The reduced model will be evaluated against the

335 full model using Akaike's Information Criterion, and by examining the effect on discrimination

336 and calibration.

337

338 DemPoRT will be developed and validated using temporal split samples, however the final

339 regression coefficients will use the full data set to maximize follow-up duration. A cohort-

340 specific intercept and/or interaction term may be included in the final model if the derivation and

341 validation cohorts differ; otherwise, the final combined model will maintain the same predictors

342 and form as the derivation model.

343

*Assessment of Predictive Performance*

345 Predictive performance in the derivation and validation cohorts will be assessed and reported

346 using overall measures of predictive accuracy, discrimination and calibration. Accuracy will be

347 assessed with Nagelkerke's $R^2$ [47] and the Brier score[48]. Discrimination will be assessed using the

348 concordance statistic. Model calibration is especially important in the development of prognostic

349 models, as probabilities of future risk are of primary interest[33,49,50]. Calibration will be assessed

350 by comparing the observed and predicted risk of dementia within vigintiles (20 groups of equal

351 frequency) of predicted risk with emphasis on visual inspection of plots rather than formal

352 statistical significance testing, which can be influenced by large sample sizes[32]. Calibration

353 slopes will be generated by regressing the outcome in the validation cohort on the predicted

354 dementia risk, reflecting the combined effect of overfitting to the derivation data as well as true

355 differences in effects of predictors. Deviation of the slope from 1 (perfect calibration) will be

356 tested using a Wald or likelihood ratio test. Calibration within predefined subgroups of

357 importance to clinicians and policy makers (e.g., age group, health behavior, sociodemographic

358 groups and health conditions) will additionally be evaluated. The clinically relevant standard of

359 calibration was defined as less than 20% difference between observed and predicted estimates

360 within subgroups with a dementia prevalence of at least 5%. All model performance measures

361 will be calculated using the first of the multiply imputed data sets.

362

363 *Model Presentation*

364 The final regression model, derived from the combined sample of the derivation and validation

365 cohorts, will be presented using estimated hazard ratios and 95% confidence intervals, along

366 with results for the derivation and validation cohorts separately. We have found, however, this

367 usual presentation less meaningful when presenting complex models[39]. To allow interpretation of

368    the estimated effect of each predictor, model behavior will additionally be described using

369    interactive visual tools to display the shape of the effect of each predictor[51]. The regression

370    formula will also be published and used as the basis for web-based implementation.

371

372    **Analyses Beyond Initial Model Development**

373    We will conduct further analyses exploring the added predictive ability of novel risk factors that

374    were ascertained in single CCHS cycles (e.g. sedentary activity, cognitive stimulation, sleep

375    quality and duration, deafness), as well as risk factors that can be ascertained through linkage of

376    additional data sources and similar cohorts (e.g. air pollution, detailed dietary consumption, lipid

377    levels, blood pressure). In addition, sensitivity analysis of the age at survey administration cutoff

378    used for cohort creation will be performed.

379

380    Once developed, DemPoRT will be used to project dementia incidence under different

381    assumptions by entering counterfactual risk factor levels in to the algorithm at the population

382    level, or at individual level and summed, and will be integrated in to POHEM for

383    microsimulation modelling of prevalence projections.

384

385    A second, causal model (DemPoRT-C) will also be created to assess the relative contribution of

386    lifestyle, socio-demographic and health factors to dementia incidence. Development will exclude

387    variables believed to be in the causal pathway of dementia occurrence (e.g., self-rated health and

388    functional measures) to reduce the attenuation of hazards from upstream risk factors, but will

389    otherwise be the same as in the predictive model. DemPoRT-C will be applied to the most recent

390    unlinked national CCHS survey.

391

**ETHICS AND MODEL DISSEMINATION**

393 The DemPoRT project advisory committee has been created to ensure that the models meet the

394 needs of knowledge users. This committee has worked with the study team to identify predictors

395 of dementia based on scientific and policy importance and will aid in the identification of

396 important target populations and the establishment of policy-relevant differences for calibration

397 studies.

398

399 DemPoRT results will be submitted for publication in peer-review journals and presented at

400 scientific meetings. A web-based individual-level calculator will be created if the models are

401 appropriate for individual use. Although DemPoRT emphasizes risk prediction at the population-

402 level, we have found that individual-level calculators are an effective engagement and translation

403 tool for both the general public and knowledge users.

404

**CONCLUSIONS**

406 To the best of our knowledge, DemPoRT will be the first population-based algorithm designed to

407 predicting and projecting dementia incidence at the population level. The DemPoRT models will

408 produce estimates of future dementia burden that we believe will be more accurate than existing

409 estimates, will assess the contribution of specific risk factors to the population risk, and identify

410 groups at high risk of developing dementia. Although a rigorous approach to model development

411 will be used, further validation will be needed to assess generalizability, and calibration will be

412 required for application in other jurisdictions.

413

414 **CONTRIBUTIONS**

415 SF drafted and revised the manuscript, and contributed to study design and protocol

416 development. NM contributed to study design, protocol development and provided

417 data/statistical support. AH, MT, DM and GH contributed to the design of the study and protocol

418 development. PT is the lead investigator of the study and was responsible for the conception of

419 the project, the grant application, study design and protocol development. All authors provided

420 critical reviews of the manuscript and approved the final version.

421

434

435 **COMPETING INTERESTS**

436 None declared.

437

**ETHICS APPROVAL**

Research ethics approval has been granted by the Ottawa Health Science Network Research

Ethics Board.

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460    **REFERNCES**

461    1.    Brookmeyer, R., Gray, S. & Kawas, C. Projections of Alzheimer's disease in the United

462          States and the public health impact of delaying disease onset. *Am. J. Public Health* **88,**

463          1337–1342 (1998).

464    2.    Brayne, C. The elephant in the room - healthy brains in later life, epidemiology and public

465          health. *Nat. Rev. Neurosci.* **8,** 233–9 (2007).

466    3.    Norton, S., Matthews, F. E. & Brayne, C. A commentary on studies presenting projections

467          of the future prevalence of dementia. *BMC Public Health* **13,** 1 (2013).

468    4.    Sloane, P. D. *et al.* The public health impact of Alzheimer's disease, 2000-2050: potential

469          implication of treatment advances. *Annu. Rev. Public Health* **23,** 213–31 (2002).

470    5.    Mura, T., Dartigues, J. F. & Berr, C. How many dementia cases in France and Europe?

471          Alternative projections and scenarios 2010-2050. *Eur. J. Neurol.* **17,** 252–259 (2010).

472    6.    Hebert, L. E., Scherr, P. A., Bienias, J. L., Bennett, D. A. & Evans, D. A. Alzheimer

473          disease in the US population: prevalence estimates using the 2000 census. *Arch. Neurol.*

474          **60,** 1119–22 (2003).

475    7.    Brookmeyer, R., Johnson, E., Ziegler-Graham, K. & Arrighi, H. M. Forecasting the global

476          burden of Alzheimer's disease. *Alzheimer's Dement.* **3,** 186–191 (2007).

477    8.    Dewey, M. E. & Chen, C.-M. Neurosis and mortality in persons aged 65 and over living in

478          the community: a systematic review of the literature. *Int. J. Geriatr. Psychiatry* **19,** 554–7

479          (2004).

480    9.    Norton, S., Matthews, F. E., Barnes, D. E., Yaffe, K. & Brayne, C. Potential for primary

481          prevention of Alzheimer's disease: An analysis of population-based data. *Lancet Neurol.*

482          **13,** 788–794 (2014).

483    10.    Joly, P. *et al.* Prevalence projections of chronic diseases and impact of public health

484           intervention. *Biometrics* **69,** 109–17 (2013).

485    11.    Lee, Y. The recent decline in prevalence of dementia in developed countries: implications

486           for prevention in the Republic of Korea. *J. Korean Med. Sci.* **29,** 913–8 (2014).

487    12.    Barnes, D. E. & Yaffe, K. The projected effect of risk factor reduction on Alzheimer's

488           disease prevalence. *Lancet. Neurol.* **10,** 819–28 (2011).

489    13.    Manuel, D. G. *et al.* Projections of preventable risks for cardiovascular disease in Canada

490           to 2021: a microsimulation modelling approach. *C. Open* **2,** E94–E101 (2014).

491    14.    Tang, E. Y. H. *et al.* Current developments in dementia risk prediction modelling: An

492           updated systematic review. *PLoS One* **10,** 1–31 (2015).

493    15.    Stephan, B. C. M., Kurth, T., Matthews, F. E., Brayne, C. & Dufouil, C. Dementia risk

494           prediction in the population: are screening models accurate? *Nat. Rev. Neurol.* **6,** 318–326

495           (2010).

496    16.    Anstey, K. J. *et al.* A self-report risk index to predict occurrence of dementia in three

497           independent cohorts of older adults: The ANU-ADRI. *PLoS One* **9,** (2014).

498    17.    Wolfsgruber, S. *et al.* The CERAD neuropsychological assessment battery total score

499           detects and predicts alzheimer disease dementia with high diagnostic accuracy. *Am. J.*

500           *Geriatr. Psychiatry* **22,** 1017–1028 (2014).

501    18.    Barnes, D. E. *et al.* Predicting risk of dementia in older adults: The late-life dementia risk

502           index. *Neurology* **73,** 173–179 (2009).

503    19.    Chary, E. *et al.* Short- versus long-term prediction of dementia among subjects with low

504           and high educational levels. *Alzheimers. Dement.* **9,** 562–71 (2013).

505    20.    Jessen, F. *et al.* Prediction of Dementia in Primary Care Patients. *PLoS One* **6,** e16852

506    (2011).

507    21.    Song, X., Mitnitski, A. & Rockwood, K. Nontraditional risk factors combine to predict

508    Alzheimer disease and dementia. *Neurology* **77,** 227–234 (2011).

509    22.    Tierney, M. C., Moineddin, R. & McDowell, I. Prediction of all-cause dementia using

510    neuropsychological tests within 10 and 5 years of diagnosis in a community-based sample.

511    *J. Alzheimer's Dis.* **22,** 1231–1240 (2010).

512    23.    Meng, X., D'Arcy, C., Morgan, D. & Mousseau, D. D. Predicting the risk of dementia

513    among Canadian seniors: a useable practice-friendly diagnostic algorithm. *Alzheimer Dis.*

514    *Assoc. Disord.* **27,** 23–9 (2013).

515    24.    Kivipelto, M. *et al.* Risk score for the prediction of dementia risk in 20 years among

516    middle aged people: a longitudinal, population-based study. *Lancet Neurol.* **5,** 735–741

517    (2006).

518    25.    Reitz, C. *et al.* A Summary Risk Score for the Prediction of Alzheimer Disease in Elderly

519    Persons. *Arch. Neurol.* **67,** 835–41 (2010).

520    26.    Walters, K. *et al.* Predicting dementia risk in primary care: development and validation of

521    the Dementia Risk Score using routinely collected data. *BMC Med.* **14,** 1–12 (2016).

522    27.    Béland, Y. Canadian community health survey--methodological overview. *Heal. reports /*

523    *Stat. Canada, Can. Cent. Heal. Inf. = Rapp. sur la sant?? / Stat. Canada, Cent. Can.*

524    *d'information sur la sant??* **13,** 9–14 (2002).

525    28.    Jaakkimainen, R. L. *et al.* Identification of Physician-Diagnosed Alzheimer's Disease and

526    Related Dementias in Population-Based Administrative Data: A Validation Study Using

527    Family Physicians' Electronic Medical Records. *J. Alzheimer's Dis.* **54,** 337–349 (2016).

528    29.    Connolly, A., Gaehl, E., Martin, H., Morris, J. & Purandare, N. Underdiagnosis of

529   dementia in primary care: Variations in the observed prevalence and comparisons to the

530   expected prevalence. *Aging Ment. Health* **15,** 978–984 (2011).

531   30.   Kosteniuk, J. G. *et al.* Incidence and prevalence of dementia in linked administrative

532         health data in Saskatchewan, Canada: a retrospective cohort study. *BMC Geriatr.* **15,** 73

533         (2015).

534   31.   Harrell, F. E. *Regression Modeling Strategies with applicaitons to linear models, logistic

535         regression and survival analysis.* (Springer, 2001).

536   32.   Vergouwe, Y., Steyerberg, E. W., Eijkemans, M. J. C. & Habbema, J. D. F. Substantial

537         effective sample sizes were required for external validation studies of predictive logistic

538         regression models. *J. Clin. Epidemiol.* **58,** 475–483 (2005).

539   33.   Steyerberg, E. W. *Clinical Prediction Models.* (Springer, 2009).

540   34.   Harrell, F. E. *Hmisc: Harrell Miscellaneous. R package version 4.0-2.* (2016).

541   35.   R Core Team. R: A language and environment for statistical computing. (2016).

542   36.   Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. M. Transparent reporting of

543         a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The

544         TRIPOD Statement. *Eur. Urol.* **67,** 1142–1151 (2015).

545   37.   Exalto, L. G. *et al.* Midlife risk score for the prediction of dementia four decades later.

546         *Alzheimer's Dement.* **10,** 562–570 (2014).

547   38.   Exalto, L. G. *et al.* Risk score for prediction of 10 year dementia risk in individuals with

548         type 2 diabetes: A cohort study. *Lancet Diabetes Endocrinol.* **1,** 183–190 (2013).

549   39.   Taljaard, M. *et al.* Cardiovascular Disease Population Risk Tool (CVDPoRT): predictive

550         algorithm for assessing CVD risk in the community setting. A study protocol. *BMJ Open*

551         **4,** e006701 (2014).

552    40.    Manuel, D. G. *et al.* Measuring Burden of Unhealthy Behaviours Using a Multivariable

553           Predictive Approach: Life Expectancy Lost in Canada Attributable to Smoking, Alcohol,

554           Physical Inactivity, and Diet. *PLoS Med.* **13,** 1–27 (2016).

555    41.    Pampalon, R., Hamel, D., Gamache, P. & Raymond, G. A deprivation index for health

556           planning in Canada. *Chronic Dis. Can.* **29,** 178–91 (2009).

557    42.    Rubin, D. B. & Schenker, N. Multiple imputation in health-care databases: an overview

558           and some applications. *Stat. Med.* **10,** 585–98 (1991).

559    43.    Therneau, T. M., Grambsch, P. M. & Fleming, T. R. Martingale-based residuals for

560           survival models. *Biometrika* **77,** 147–160 (1990).

561    44.    Allison, P. D. *Survival analysis using SAS : A practical guide*. (SAS Institute, 2010).

562    45.    Van Houwelingen, J. C. & Le, C. S. Predictive value of statistical models. *Stat Med* **9,**

563           1303–1325 (1990).

564    46.    Ambler, G., Brady, A. R. & Royston, P. Simplifying a prognostic model: A simulation

565           study based on clinical data. *Stat. Med.* **21,** 3803–3822 (2002).

566    47.    Nagelkerke, N. J. D. A Note on a General Definition of the Coefficient of Determination.

567           *Biometrika* **78,** 691–692 (1991).

568    48.    Arkes, H. R. *et al.* The covariance decomposition of the probability score and its use in

569           evaluating prognostic estimates. SUPPORT Investigators. *Med. Decis. Making* **15,** 120–31

570           (1995).

571    49.    Cook, N. R. Statistical evaluation of prognostic versus diagnostic models: beyond the

572           ROC curve. *Clin. Chem.* **54,** 17–23 (2008).

573    50.    Cook, N. R. Comment: Measures to summarize and compare the predictive capacity of

574           markers. *Int. J. Biostat.* **6,** Article 22; discussion Article 25 (2010).

575    51.    Krause, J., Perer, A. & Bertini, E. Using Visual Analytics to Interpret Predictive Machine

576           Learning Models. *arXiv Prepr. arXiv1606.05685.* (2016).

577

578

579

580

581

582

583

584

585

586

587

588

589

590

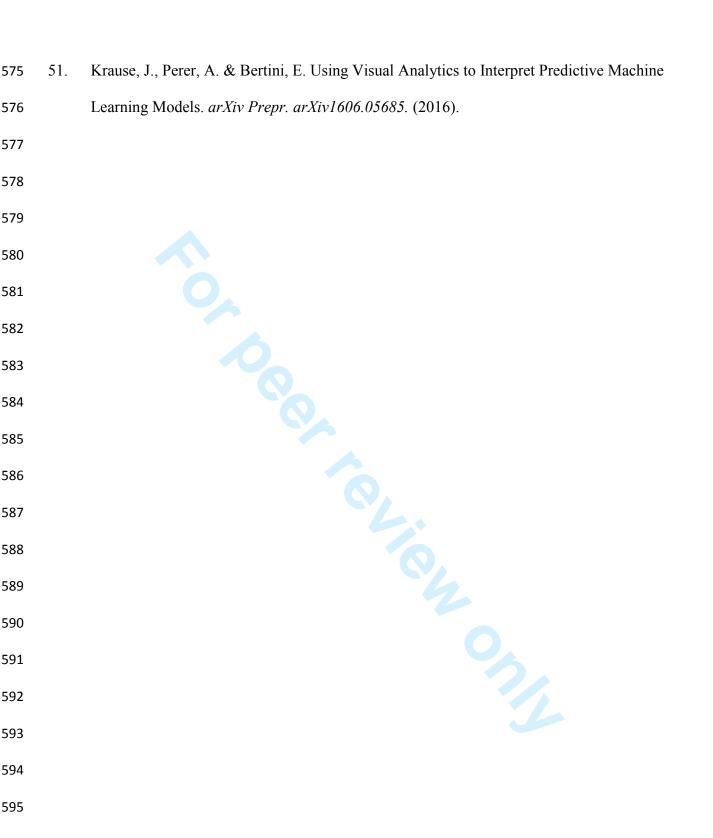591

592

593

594

595

596

597

598

599    Table 1. Pre-specification of predictor variables for DemPoRT with initial degrees of freedom (df) allocation

| Variable | Scale | Initial Variable Specification | df |
|---|---|---|---|
| **Socio-demographic Factors** | | | |
| Age | Continuous | **5 knot spline:** Valid range: 55-102 (male), 55-101 (female) | 4 |
| Sex | Categorical | **Stratified:** Male; Female | NA |
| Ethnicity | Categorical | **7 categories:** Caucasian; African-American; Chinese; Aboriginal; Japanese/Korean/South East Asian/Filipino; Other/Multiple origin/Unknown/Latin American; South Asian/Arab/West Asian | 6 |
| Immigrant | Dichotomous | Yes; No | 1 |
| Education | Categorical | **4 categories:** Less than secondary school; Secondary school graduation; Some postsecondary; Postsecondary graduation | 3 |
| Marital Status | Categorical | **4 categories:** Now married/Common-law; Separated/Divorced; Widowed; Single | 3 |
| Neighborhood Social and Material Deprivation[41] | Ordinal | **3 categories:** Low (1st or 2nd quintile); High 4th or 5th quintile; Moderate (3rd quintile) | 2 |
| **General Health** | | | |
| Sense of belonging to local community | Ordinal | **4 categories:** Very strong; Somewhat strong; Somewhat weak; Very weak | 3 |
| Self-perceived stress | Ordinal | **5 categories:** Not at all stressful; Not very stressful; A bit stressful; Quite a bit stressful; Extremely stressful | 4 |
| Self-rated health | Ordinal | **5 categories:** Poor; Fair; Good; Very Good; Excellent | 4 |
| **Health Behaviors** | | | |
| Pack years of smoking | Continuous | **3 knot spline:** Valid range: 0-112 (male), 0-78 (female) | 2 |
| Smoking status | Categorical | **4 categories:** Non-smoker; Current smoker; Former smoker quit <5 years ago; Former smoker quit >5 years ago | 3 |
| Alcohol consumption (number of drinks last week) | Continuous | **3 knot spline:** Valid range: 0-50 (male), 0-24 (female) | 2 |
| Former drinker | Dichotomous | Yes; No | 1 |
| Consumption of fruit, salad, carrot and other vegetables (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-48 (male), 0-31 (female) | 2 |
| Potato consumption (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-2 | 2 |
| Juice consumption (average daily consumption | Continuous | **3 knot spline:** Valid range: 0-6 (male), 0-5 (female) | 2 |
| Leisure physical activity (average daily METs (kcal/kg/day)) | Continuous | **3 knot spline:** Valid range: 0-16 (male), 0-12 (female) | 2 |
| **Functional Measures** | | | |
| Personal hygiene and care | Dichotomous | Does not need help; Needs help | 1 |
| Locomotion in the home | Dichotomous | Does not need help; Needs help | 1 |
| Meal preparation | Dichotomous | Does not need help; Needs help | 1 |
| Running errands | Dichotomous | Does not need help; Needs help | 1 |
| Ordinary housework | Dichotomous | Does not need help; Needs help | 1 |
| Heavy housework | Dichotomous | Does not need help; Needs help | 1 |
| Finances | Dichotomous | Does not need help; Needs help | 1 |
| **Health Conditions** | | | |
| Heart disease | Dichotomous | Yes; No | 1 |
| Stroke | Dichotomous | Yes; No | 1 |
| Diabetes | Dichotomous | Yes; No | 1 |
| Mood disorder | Dichotomous | Yes; No | 1 |
| High blood pressure | Dichotomous | Yes; No | 1 |
| Body mass index | Continuous | **3 knot spline:** Valid range: 10-44 (male), 10-47 (female) | 2 |
| **Design** | | | |
| Survey year | Ordinal | **4 categories:** 2000/01, 2002/03, 2004/05, 2006/07 | 3 |

600    DemPoRT, Dementia Population Risk Tool; df, degrees of freedom; MET, metabolic equivalent task

## TRIPOD Checklist: Prediction Model Development and Validation

| Section/Topic | Item | | Checklist Item | Page |
|---|---|---|---|---|
| **Title and abstract** | | | | |
| Title | 1 | D;V | Identify the study as developing and/or validating a multivariable prediction model, the target population, and the outcome to be predicted. | 3 |
| Abstract | 2 | D;V | Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions. | 3 |
| **Introduction** | | | | |
| Background and objectives | 3a | D;V | Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models. | 5-7 |
| | 3b | D;V | Specify the objectives, including whether the study describes the development or validation of the model or both. | 7,8 |
| **Methods** | | | | |
| Source of data | 4a | D;V | Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable. | 8,9 |
| | 4b | D;V | Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up. | 8-10 |
| Participants | 5a | D;V | Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres. | 8,9 |
| | 5b | D;V | Describe eligibility criteria for participants. | 8-10 |
| | 5c | D;V | Give details of treatments received, if relevant. | NA |
| Outcome | 6a | D;V | Clearly define the outcome that is predicted by the prediction model, including how and when assessed. | 9,10 |
| | 6b | D;V | Report any actions to blind assessment of the outcome to be predicted. | 12 |
| Predictors | 7a | D;V | Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured. | 11,12, 25 |
| | 7b | D;V | Report any actions to blind assessment of predictors for the outcome and other predictors. | 12 |
| Sample size | 8 | D;V | Explain how the study size was arrived at. | 8-12 |
| Missing data | 9 | D;V | Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method. | 12,13 |
| Statistical analysis methods | 10a | D | Describe how predictors were handled in the analyses. | 13,14, 25 |
| | 10b | D | Specify type of model, all model-building procedures (including any predictor selection), and method for internal validation. | 14,15 |
| | 10c | V | For validation, describe how the predictions were calculated. | 16 |
| | 10d | D;V | Specify all measures used to assess model performance and, if relevant, to compare multiple models. | 16 |
| | 10e | V | Describe any model updating (e.g., recalibration) arising from the validation, if done. | NA |
| Risk groups | 11 | D;V | Provide details on how risk groups were created, if done. | 16 |
| Development vs. validation | 12 | V | For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors. | 16,17 |
| **Results** | | | | |
| Participants | 13a | D;V | Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful. | 10 |
| | 13b | D;V | Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome. | NA |
| | 13c | V | For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome). | NA |
| Model development | 14a | D | Specify the number of participants and outcome events in each analysis. | NA |
| | 14b | D | If done, report the unadjusted association between each candidate predictor and outcome. | NA |
| Model specification | 15a | D | Present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point). | NA |
| | 15b | D | Explain how to the use the prediction model. | NA |
| Model performance | 16 | D;V | Report performance measures (with CIs) for the prediction model. | NA |
| Model-updating | 17 | V | If done, report the results from any model updating (i.e., model specification, model performance). | NA |
| **Discussion** | | | | |
| Limitations | 18 | D;V | Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data). | 4,17 |
| Interpretation | 19a | V | For validation, discuss the results with reference to performance in the development data, and any other validation data. | NA |
| | 19b | D;V | Give an overall interpretation of the results, considering objectives, limitations, results from similar studies, and other relevant evidence. | NA |
| Implications | 20 | D;V | Discuss the potential clinical use of the model and implications for future research. | 5-7 |
| **Other information** | | | | |
| Supplementary information | 21 | D;V | Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets. | 4,18 |
| Funding | 22 | D;V | Give the source of funding and the role of the funders for the present study. | 19 |

*Items relevant only to the development of a prediction model are denoted by D, items relating solely to a validation of a prediction model are denoted by V, and items relating to both are denoted D;V. We recommend using the TRIPOD Checklist in conjunction with the TRIPOD Explanation and Elaboration document.

**BMJ Open**

# Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm Assessing Dementia Risk in the Community

SCHOLARONE™
Manuscripts

1    **Dementia Population Risk Tool (DemPoRT): Study Protocol for a Predictive Algorithm**

2    **Assessing Dementia Risk in the Community**

3

4    Stacey Fisher, MSc[1,2,3] stacey.fisher@uottawa.ca

5    Amy Hsu, PhD[1,2] ahsu@toh.ca

6    Nassim Mojaverian, MSc[2] namojaverian@ohri.ca

7    Monica Taljaard, PhD[1,3] mtaljaard@ohri.ca

8    Gregory Huyer, PhD[4] ghuye047@uottawa.ca

9    Doug Manuel, MD[1,2,3,5,6] dmanuel@ohri.ca

10    Peter Tanuseputro, MD[1,2,6,7,8] ptanuseputro@ohri.ca

11

12    [1] Ottawa Hospital Research Institute, Ottawa, Ontario, Canada

13    [2] Institute for Clinical Evaluative Sciences, Ottawa, Ontario, Canada

14    [3] Department of Epidemiology and Community Medicine, University of Ottawa, Ottawa,

15    Ontario, Canada

16    [4] Telfer School of Management, University of Ottawa, Ottawa, Ontario, Canada

17    [5] Statistics Canada, Ottawa, Ontario, Canada

18    [6] Department of Family Medicine, University of Ottawa, Ottawa, Ontario, Canada

19    [7] Department of Medicine, University of Ottawa, Ottawa, Ontario, Canada

20    [8] Bruyère Research Institute, Ottawa, Ontario, Canada

21

22    Corresponding author:

23    Dr. Peter Tanuseputro

24    1053 Carling Ave.

25    Box 693

26    Ottawa, ON K1Y 4E9, Canada

27    ptanuseputro@ohri.ca

28    Phone: 613-798-5555

29

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47    **ABSTRACT**

48    **Introduction:** The burden of disease from dementia is a growing global concern as incidence

49    increases dramatically with age and average life expectancy has been increasing around the

50    world. Planning for an aging population requires reliable projections of dementia prevalence;

51    however, existing population projections are simple and have poor predictive accuracy. The

52    Dementia Population Risk Tool (DemPoRT) will predict incidence of dementia in the population

53    setting using multivariable modeling techniques, and will be used to project dementia

54    prevalence.

55    **Methods and Analysis:** The derivation cohort will consist of elderly Ontario respondents of the

56    Canadian Community Health Survey (CCHS) (2001, 2003, 2005, 2007; 18 764 males and 25 288

57    females). Pre-specified predictors include sociodemographic, general health, behavioral,

58    functional and health condition variables. Incident dementia will be identified through individual

59    linkage of survey respondents to population-level administrative health care databases (1 797 and

60    3 281 events, and 117 795 and 166 573 person-years of follow-up, for males and females,

61    respectively until March 31, 2014). Using time of first dementia capture as the primary outcome

62    and death as a competing risk, sex-specific proportional hazards regression models will be

63    estimated. The 2008/2009 CCHS survey will be used for validation (approximately 4 600 males

64    and 6 300 females). Overall calibration and discrimination will be assessed as well as calibration

65    within predefined subgroups of importance to clinicians and policy makers.

66    **Ethics and Dissemination:** Research ethics approval has been granted by the Ottawa Health

67    Science Network Research Ethics Board. DemPoRT results will be submitted for publication in

68    peer-review journals and presented at scientific meetings. The algorithm will be assessable

69    online for both population and individual uses.

70    **Trial Registration Number:** ClinicalTrials.gov NCT03155815.

71

72    **STRENGTHS AND LIMITATIONS**

73    -    The Dementia Population Risk Tool (DemPoRT) will be developed and validated using

74         predictors from large population-based community health surveys that are individually

75         linked to routinely-collected health administration data in Ontario. To our knowledge,

76         DemPoRT will be the first algorithm designed to predict and project dementia incidence

77         at the population-level.

78    -    Although repeat predictor assessment and detailed cognitive testing to ascertain dementia

79         diagnoses is preferable, it is not available or feasible at the population level.

80    -    Statistical overfitting is a concern, however full pre-specification of the analysis plan and

81         predictors will limit this risk.

82    -    Although a rigorous approach to model development will be used, further validation will

83         be needed to assess generalizability, and calibration will be required for application in

84         other jurisdictions.

85    -    DemPoRT will be used to produce improved estimates of future dementia burden, will

86         assess the contribution of specific risk factors to the population risk, and will identify

87         population subgroups at high risk of developing dementia. This information will be used

88         by policymakers to prepare for and reduce dementia impact.

89

90

91

92

93    **INTRODUCTION**

94    The burden of disease from dementia is a growing global concern as incidence increases

95    dramatically with age and average life expectancy has been increasing around the world[1,2].

96    Planning for an aging population requires reliable projections of dementia burden and the

97    implications for resource requirements. Existing population-level projections for dementia,

98    however, are overly simplistic and likely inaccurate[3].

99

100    **Limitations of Current Dementia Projection Methodology**

101    Almost all existing dementia projections have used extrapolation and macrosimulation methods,

102    which are simplistic and make assumptions that may not hold true into the future[3]. Most

103    extrapolations simply apply current age- and sex-specific prevalence estimates of dementia to

104    future population projections. Macrosimulations typically use estimates of dementia incidence

105    and mortality, stratified by age and sex, to simulate disease prevalence as the population ages[1,4–

106    6]. Projections from extrapolations incorrectly assume that the risk of mortality among those with

107    and without dementia are equivalent[7,8], and both methods assume that the age and sex-specific

108    prevalence of dementia risk factors will not change with time. The assumption of stable risk

109    factor prevalence is widely thought to be the major source of error in existing dementia

110    projections[3,9–11].

111

112    Changing trends of dementia risk factors has the potential to have a dramatic impact on dementia

113    prevalence estimates, as up to 50% of dementia cases have been attributed to modifiable

114    factors[9,12], and the prevalence of several factors has been projected to change significantly in the

115    near future. For example, the population prevalence of diabetes and obesity in Canada has been

116  projected to increase, while smoking, hypertension and dyslipidemia have been projected to

117  decline[13]. Consideration of risk factor prevalence is therefore important to improve the accuracy

118  of dementia projections, and simple extrapolations and macrosimulations are often inadequate.

119

**Predictive Multivariable Modeling of Dementia Incidence**

121  Population-based **predictive risk algorithms** examine the effect of risk factors on dementia

122  incidence, and can be used for dementia burden projection. Population-based data that contain

123  detailed risk factor information, such as health surveys, are linked at the individual-level to

124  administrative data that capture dementia development. A multivariable model of dementia

125  incidence is derived, validated against external data, and predictive performance is assessed.

126  Once developed, the algorithm can be used to project disease incidence and prevalence. To

127  obtain prevalence projections, the algorithm can be integrated in to a microsimulation model

128  such as Statistics Canada's Population Health Models (POHEM). POHEM dynamically models

129  individual life trajectories of a population representative of Canada including births, deaths and

130  migration, disease incidence and progression, and exposure to risk factors, facilitating detailed

131  examination of the influence of changing risk factor prevalence on future dementia prevalence.

132

133  Predictive risk algorithms can also be used to describe the risk of dementia in the population,

134  assess the contribution of specific risk factors to the population risk, identify high-risk groups,

135  and evaluate risk reduction strategies.

136

**Existing Dementia Prediction Models**

138    Many models have been developed to predict risk of dementia[14–26], most with the primary goal of

139    identifying individuals in the clinical setting at high risk. They have varying discriminative

140    ability (c-statistics ranging from 0.49[16] to 0.89[17]) and have generally been derived from small

141    samples, rarely including more than a few thousand individuals. Existing models are therefore

142    simplistic, including few predictors and rarely including interaction or non-linear terms Existing

143    models thus facilitates understanding and use by physicians in clinical practice, but limits

144    discriminatory ability and predictive accuracy. Walters et al[26] developed an algorithm for

145    predicting 5-year dementia risk among individuals 60-79 years of age in the United Kingdom

146    using an enormous derivation dataset of 800 000 individuals, and a simple model. The derivation

147    model had a c-statistic of 0.84 (95% CI: 0.81, 0.87), but a low positive predictive value at most

148    risk thresholds, and therefore is poor at identifying those at high risk of dementia. Additionally,

149    as most dementia risk models are intended for use in the clinical setting, many include results

150    from neuropsycological tests[17–23], MRI findings[18] and APOE genotype[18,24,25]. The inclusion of

151    these variables, however, limits the application of these models as these variables are not

152    available at the population-level.

153

154    The objective of this study is to develop and validate the Dementia Population Risk Tool

155    (DemPoRT) algorithm to predict dementia incidence in the population setting. This will be done

156    using multivariable modelling techniques, linking self-reported risk factors captured by a

157    population-based health survey in Canada with administrative databases across healthcare sectors

158    that capture healthcare diagnosed dementia. DemPoRT will be developed with a using a large

159    population-based dataset using only variables that are available at the population-level, allowing

160    for population-level application. DemPoRT will also utilize many methodological improvements

161    over existing models. This protocol pre-specifies the predictor variables and analytic plan for

162    model development, reducing the potential for overfitting and bias, and improving transparency.

163    Interaction terms and flexible functions for continuous predictors will be investigated, increasing

164    potential discriminative ability. The pre-specified analytic plan avoids data-driven variable

165    selection procedures, further reducing the potential for bias.

166

167    To our knowledge, the DemPoRT predictive model will be the first algorithm designed to predict

168    and project dementia incidence at the population-level. It will be used to estimate the future

169    burden of dementia using techniques that consider changes in risk factor prevalence and will

170    identify modifiable risk factors that can be targeted by individuals, clinicians and policy makers

171    to reduce the burden of dementia.

172

173    **METHODS AND ANALYSIS**

174    **Study Design**

175    Two DemPoRT models, one for males and females, will be derived and validated using

176    population-based data in Ontario, Canada, a multicultural province with 13.6 million residents.

177    Predictors will be obtained from the Canadian Community Health Surveys (CCHS), and

178    outcomes (i.e., diagnosis of dementia) will be obtained from routinely-collected health care data.

179

180    The derivation cohorts will consist of eligible respondents of the 2001, 2003, 2005 and 2007

181    CCHS (Cycles 1.1, 2.1, 3.1 and 4.1), while validation cohorts will consist of respondents to the

182    2008/2009 cycle. The CCHS is a national, cross-sectional survey developed by Statistics Canada

183    to collect information related to health and health care utilization of the Canadian population.

184   The survey has a multistage stratified cluster design that represents approximately 98% of the

185   Canadian population aged 12 years and over and attained an average response rate of 79% over

186   the study period. The CCHS is conducted through telephone and in-person interviews, and all

187   responses are self-reported. The details of survey methodology have been published elsewhere[27].

188   Survey respondents will be excluded if they are less than 55 years of age at survey

189   administration, self-reported a history of dementia, or are not eligible for Ontario's universal

190   health insurance. If a respondent was included in more than one CCHS cycle, only their earliest

191   survey response will be used.

192

193   **Outcome**

194   Survey respondents diagnosed with dementia will be identified through individual linkage to

195   several population-based administrative databases at the Institute for Clinical Evaluative

196   Sciences (ICES). Dementia case ascertainment is based on a validated definition: 1 hospital

197   record OR 3 physician claim records at least 30 days apart within a 2-year period OR a

198   dispensing record for a cholinesterase inhibitor from Ontario Drug Benefit (ODB). This

199   definition has a sensitivity of 79.3% and a specificity of 99.1% when validated against

200   emergency medical record (EMR) data[28]. Due to known underdiagnosis of dementia[29,30], we will

201   supplement this definition by adding survey respondents with dementia codes captured on home

202   care and long-term care assessments (dementia flag AND Cognitive Performance Scale [CPS]

203   score $\geq 2$) using the Resident Assessment Instrument-Home Care (RAI-HC) database and the

204   Continuing Care Reporting System (CCRS), respectively. We have found this addition adds

205   substantially (approximately 18%) to the number of dementia cases captured.

206

207   Survey respondents with dementia will be excluded if they meet the criteria for dementia within

208   two years of survey administration (to remove potentially prevalent cases) or are younger than 65

209   years of age at the time of dementia diagnosis (to exclude early onset dementia which likely has

210   a different set of risk factors). Eligible survey respondents will be followed from the date of

211   survey administration or age 65, whichever came later, until the earliest date of: dementia

212   ascertainment, death (defined as competing risk), loss to follow-up (defined as loss of healthcare

213   eligibility) or end of study (March 31, 2014).

214

215   **Sample Size**

216   The male and female derivation cohorts consist of 18 764 and 25 288 respondents, and 117 795

217   and 166 573 person-years of follow-up, respectively. For predictive models with time to event

218   outcomes the number of participants experiencing the event should exceed 10 times the number

219   of degrees of freedom to ensure adequate sample size[31]. The number of dementia events in the

220   derivation cohort is 1 797 for men and 3 281 for women; therefore, the maximum number of

221   total degrees of freedom for each of the DemPoRT models is 179 and 328, respectively, which

222   we do not anticipate surpassing.

223

224   The validation cohorts will consist of approximately 4 600 males and 6 300 females, and 15 000

225   and 21 000 person-years of follow-up, respectively. Vergouwe *et al*[32] recommend a minimum of

226   100 events and 100 non-events for external validation studies. We expect approximately 225

227   events for men and 400 for women in our validation cohort.

228

229   **Analysis Plan**

230  The analysis plan was developed following guidelines by Harrell[31] and Steyerberg[33] after

231  accessing the derivation data set, but prior to model fitting or descriptive analyses involving

232  exposure-outcome associations. This was done to avoid Type 1 error introduced by data-driven

233  variable selection or model specification. Key considerations of our analysis approach include

234  full pre-specification of the predictor variables, use of flexible functions for continuous

235  predictors, and preserving statistical properties by avoiding data-driven variable selection

236  procedures. Analysis will be conducted using Harrell's Hmisc[34] package of functions in R[35] as

237  well as SAS v9.4.

238

239  This study protocol and the reporting of our model estimation results will be guided by the

240  TRIPOD statement for multivariable predictive models[36].

241

242  *Identification of Predictors*

243  Predictor variables were identified through review of existing predictive algorithms for

244  dementia[9,16,18–22,24–26,37,38] and comparison to available data collected in the CCHS. Variable

245  inclusion was informed by consultation with subject-matter experts and the project's advisory

246  team, and informed by our previous work developing predictive models for cardiovascular

247  disease and life expectancy[39,40].

248

249  Variables with narrow distributions or insufficient variation were excluded. Obvious cases of

250  redundancy (e.g. alternate definitions of the same underlying behavior) were not included. A

251  total of 32 predictor variables were identified: 7 sociodemographic, 3 general health, 9

252  behavioral, 7 functional, 5 health conditions and 1 design variable (CCHS survey cycle). As the

253  effect of dementia risk factors varies by sex, separate models will be derived for men and

254  women. Education, rather than individual income, was selected as a predictor due to several

255  concerns with income including lack of generalizability, measurement error, stability over time

256  and substantial missing values. Neighborhood social and material deprivation is captured using

257  Pampalon's deprivation index[41]. Indicator variables for smoking status were created to allow the

258  inclusion of smoking pack-years as a continuous predictor. The models will additionally include

259  age interactions with the behavioral, functional and health condition variables as the effect of

260  these risk factors on dementia are expected to vary with age. Detailed definitions and

261  measurement of the predictor variables are presented in Table 1.

262

263  *Data Cleaning and Coding of Predictors*

264  Continuous variables will be inspected using boxplots and descriptive statistics to determine

265  values outside a plausible range. Values that are clearly erroneous will be corrected, where

266  possible, or set to missing. Continuous predictors with highly skewed distributions will be

267  truncated to the 99.5th percentile. Categorization of continuous variables will be avoided to

268  minimize the loss of predictive information. All data cleaning and coding will occur prior to

269  examining exposure-outcome associations.

270

271  *Missing Data*

272  As traditional complete cases analyses suffer from inefficiency, selection bias, and other

273  limitations[33], multiple imputation methods will be used to impute missing values using the

274  'aregImpute' function in the HMisc library[34]. This function simultaneously imputes missing

275  values using predictive mean matching and uses bootstrapping to take all aspects of uncertainty

276  in to account. The imputation model will consist of the full list of predictor variables, time to

277  event and censoring variables, as well as auxiliary variables that are not predictors, but may

278  nevertheless be useful in generating imputed values (e.g., income). The final model will be

279  estimated in each of five multiple imputation data sets and the results combined using the rules

280  developed by Rubin and Schenker[42] to account for imputation uncertainty.

281

282  *Model Specification*

283  Initial sex-specific main effects models will be fit using the pre-specified predictors and an initial

284  degree of freedom allocation for each predictor (Table 1). Decisions on initial degree of freedom

285  allocations were informed by the anticipated importance of each predictor and known dose-

286  response relationships with dementia. Continuous predictors will be flexibly modelled using

287  restricted cubic splines, with the knots placed at fixed quantiles of the distribution (e.g., 5th,

288  27.5th, 50th, 72.5th, and 95th centiles). Frequency distributions for categorical predictors will be

289  examined and categories with small numbers of respondents will be combined, with analysts

290  blinded to the number of events per category, to avoid instability in the regression analyses.

291  Ordinal variables will be specified as either linear terms or as categorical if the expected

292  association is more complex. Interactions will be restricted to linear terms. The initial model

293  specification, presented in Table 1, includes a total of 86 degrees of freedom (63 main, 23

294  interaction).

295

296  Partial association chi-square statistics for each predictor minus their degrees of freedom (to

297  level the playing field among predictors with varying degrees of freedom) will be plotted in

298  descending order. Variables with higher predictive potential will retain their initial degrees of

299   freedom, while predictors with lower predictive potential will be modeled as simple linear terms

300   or recoded by combining infrequent categories. This process of model specification does not

301   increase the Type I error rate because all predictors will be retained in the full model regardless

302   of their strength of association[31].

303

304   *Model Estimation*

305   The initial models will be estimated using competing risk Cox proportional hazards regression

306   with time to dementia ascertainment as the outcome and death as a competing risk. Alternative

307   model specifications, including subdistribution hazard and flexible parametric models, will be

308   considered. All predictors will be centered about their means. A formal check of

309   multicollinearity will be carried out using a variable clustering algorithm[31].

310

311   Proportional hazards models assume that the relative risk of the outcome between strata of

312   predictors and the baseline risk must be constant over time. Violation of this assumption has

313   been shown to produce biased results[43] although it has also been argued that the estimated

314   coefficients of time-varying variables can simply be interpreted as an average rather than

315   instantaneous hazard[44]. Plots of raw and smoothed scaled Schoenfeld residuals versus time for

316   each predictor will be assessed to test this assumption and identify non-proportionality. If a

317   violation of this assumption is identified we will consider addition of interaction terms between

318   the predictor and log-transformed time.

319

320   Although the risk of overfitting will be minimal due to pre-specification of the models and a

321   large sample size, the need for overfitting adjustment will be assessed. The degree of overfitting

322  will be estimated using the heuristic shrinkage estimator, based on the log likelihood ratio chi-

323  square statistic for the full model[45]. If shrinkage is <0.90, models will be adjusted for overfitting.

324

*Estimation of the Reduced Models*

326  Model pre-specification has advantages in limiting overfitting and spurious statistical

327  significance but can result in a final model that is overly complex, difficult to interpret, and

328  difficult to apply. Unnecessary predictor variables also distort the estimated effects of other

329  predictors making the model more computationally intensive. It is suggested that a more

330  parsimonious model that retains most of the prognostic information and performs as well as or

331  better than the full model can be derived without increasing the Type 1 error rate[31,46]. We will

332  identify a more parsimonious model using a stepdown procedure described by Ambler[46], which

333  involves deleting the variable that results in the smallest decrease in model $R^2$ until removal

334  leads to an $R^2$ that is less than a desired level. The reduced model will be evaluated against the

335  full model using Akaike's Information Criterion, and by examining the effect on discrimination

336  and calibration.

337

338  DemPoRT will be developed and validated using temporal split samples, however the final

339  regression coefficients will use the full data set to maximize follow-up duration.  A cohort-

340  specific intercept and/or interaction term may be included in the final model if the derivation and

341  validation cohorts differ; otherwise, the final combined model will maintain the same predictors

342  and form as the derivation model.

343

*Assessment of Predictive Performance*

345  Predictive performance in the derivation and validation cohorts will be assessed and reported

346  using overall measures of predictive accuracy, discrimination and calibration. Accuracy will be

347  assessed with Nagelkerke's $R^2$ [47] and the Brier score[48]. Discrimination will be assessed using the

348  concordance statistic. Model calibration is especially important in the development of prognostic

349  models, as probabilities of future risk are of primary interest[33,49,50]. Calibration will be assessed

350  by comparing the observed and predicted risk of dementia within vigintiles (20 groups of equal

351  frequency) of predicted risk with emphasis on visual inspection of plots rather than formal

352  statistical significance testing, which can be influenced by large sample sizes[32]. Calibration

353  slopes will be generated by regressing the outcome in the validation cohort on the predicted

354  dementia risk, reflecting the combined effect of overfitting to the derivation data as well as true

355  differences in effects of predictors. Deviation of the slope from 1 (perfect calibration) will be

356  tested using a Wald or likelihood ratio test. Calibration within predefined subgroups of

357  importance to clinicians and policy makers (e.g., age group, health behavior, sociodemographic

358  groups and health conditions) will additionally be evaluated. The clinically relevant standard of

359  calibration was defined as less than 20% difference between observed and predicted estimates

360  within subgroups with a dementia prevalence of at least 5%. All model performance measures

361  will be calculated using the first of the multiply imputed data sets.

362

363  *Model Presentation*

364  The final regression model, derived from the combined sample of the derivation and validation

365  cohorts, will be presented using estimated hazard ratios and 95% confidence intervals, along

366  with results for the derivation and validation cohorts separately. We have found, however, this

367  usual presentation less meaningful when presenting complex models[39]. To allow interpretation of

368    the estimated effect of each predictor, model behavior will additionally be described using

369    interactive visual tools to display the shape of the effect of each predictor[51]. The regression

370    formula will also be published and used as the basis for web-based implementation.

371

372    **Analyses Beyond Initial Model Development**

373    We will conduct further analyses exploring the added predictive ability of novel risk factors that

374    were ascertained in single CCHS cycles (e.g. sedentary activity, cognitive stimulation, sleep

375    quality and duration, deafness), as well as risk factors that can be ascertained through linkage of

376    additional data sources and similar cohorts (e.g. air pollution, detailed dietary consumption, lipid

377    levels, blood pressure). In addition, sensitivity analysis of the age at survey administration cutoff

378    used for cohort creation will be performed.

379

380    Once developed, DemPoRT will be used to project dementia incidence under different

381    assumptions by entering counterfactual risk factor levels in to the algorithm at the population

382    level, or at individual level and summed, and will be integrated in to POHEM for

383    microsimulation modelling of prevalence projections.

384

385    A second, causal model (DemPoRT-C) will also be created to assess the relative contribution of

386    lifestyle, socio-demographic and health factors to dementia incidence. Development will exclude

387    variables believed to be in the causal pathway of dementia occurrence (e.g., self-rated health and

388    functional measures) to reduce the attenuation of hazards from upstream risk factors, but will

389    otherwise be the same as in the predictive model. DemPoRT-C will be applied to the most recent

390    unlinked national CCHS survey.

**LIMITATIONS**

One of the limitations of this study will be the potential for misclassification error resulting from the use of self-reported predictors captured at one point in time and administrative data for outcome ascertainment. However, discriminating and well-calibrated algorithms have been developed using self-report information and although detailed cognitive testing to ascertain dementia diagnoses is preferable over the use of administrative data, it is not available or feasible at the population level. Another concern common to the development of highly complex risk algorithms such as DemPoRT, is the potential for statistical overfitting and increased Type 1 error, which can occur when the relationship between a predictor and the outcome influences whether it is used, and how it is fit. This risk is reduced by pre-specification of the predictors and analytic plan, as we have done in this protocol. The model will also be adjusted for overfitting if necessary, as specified previously. Lastly, although a rigorous approach to model development will be used, further validation will be needed to assess generalizability, and calibration will be required for application in other jurisdictions.

**ETHICS AND MODEL DISSEMINATION**

The DemPoRT project advisory committee has been created to ensure that the models meet the needs of knowledge users. This committee has worked with the study team to identify predictors of dementia based on scientific and policy importance and will aid in the identification of important target populations and the establishment of policy-relevant differences for calibration studies.

413   DemPoRT results will be submitted for publication in peer-review journals and presented at

414   scientific meetings. A web-based individual-level calculator will be created if the models are

415   appropriate for individual use. Although DemPoRT emphasizes risk prediction at the population-

416   level, we have found that individual-level calculators are an effective engagement and translation

417   tool for both the general public and knowledge users.

418

419   **CONCLUSIONS**

420   To the best of our knowledge, DemPoRT will be the first population-based algorithm designed to

421   predicting and projecting dementia incidence at the population level. The DemPoRT models will

422   produce estimates of future dementia burden that we believe will be more accurate than existing

423   estimates, will assess the contribution of specific risk factors to the population risk, and identify

424   groups at high risk of developing dementia. Although a rigorous approach to model development

425   will be used, further validation will be needed to assess generalizability, and calibration will be

426   required for application in other jurisdictions.

427

428   **CONTRIBUTIONS**

429   SF drafted and revised the manuscript, and contributed to study design and protocol

430   development. NM contributed to study design, protocol development and provided

431   data/statistical support. AH, MT, DM and GH contributed to the design of the study and protocol

432   development. PT is the lead investigator of the study and was responsible for the conception of

433   the project, the grant application, study design and protocol development. All authors provided

434   critical reviews of the manuscript and approved the final version.

435

448

449 **COMPETING INTERESTS**

450 None declared.

451

452 **ETHICS APPROVAL**

453 Research ethics approval has been granted by the Ottawa Health Science Network Research

454 Ethics Board.

455

456

457

458

459    **REFERNCES**

460    1.    Brookmeyer, R., Gray, S. & Kawas, C. Projections of Alzheimer's disease in the United

461          States and the public health impact of delaying disease onset. *Am. J. Public Health* **88,**

462          1337–1342 (1998).

463    2.    Brayne, C. The elephant in the room - healthy brains in later life, epidemiology and public

464          health. *Nat. Rev. Neurosci.* **8,** 233–9 (2007).

465    3.    Norton, S., Matthews, F. E. & Brayne, C. A commentary on studies presenting projections

466          of the future prevalence of dementia. *BMC Public Health* **13,** 1 (2013).

467    4.    Sloane, P. D. *et al.* The public health impact of Alzheimer's disease, 2000-2050: potential

468          implication of treatment advances. *Annu. Rev. Public Health* **23,** 213–31 (2002).

469    5.    Mura, T., Dartigues, J. F. & Berr, C. How many dementia cases in France and Europe?

470          Alternative projections and scenarios 2010-2050. *Eur. J. Neurol.* **17,** 252–259 (2010).

471    6.    Hebert, L. E., Scherr, P. A., Bienias, J. L., Bennett, D. A. & Evans, D. A. Alzheimer

472          disease in the US population: prevalence estimates using the 2000 census. *Arch. Neurol.*

473          **60,** 1119–22 (2003).

474    7.    Brookmeyer, R., Johnson, E., Ziegler-Graham, K. & Arrighi, H. M. Forecasting the global

475          burden of Alzheimer's disease. *Alzheimer's Dement.* **3,** 186–191 (2007).

476    8.    Dewey, M. E. & Chen, C.-M. Neurosis and mortality in persons aged 65 and over living in

477          the community: a systematic review of the literature. *Int. J. Geriatr. Psychiatry* **19,** 554–7

478          (2004).

479    9.    Norton, S., Matthews, F. E., Barnes, D. E., Yaffe, K. & Brayne, C. Potential for primary

480          prevention of Alzheimer's disease: An analysis of population-based data. *Lancet Neurol.*
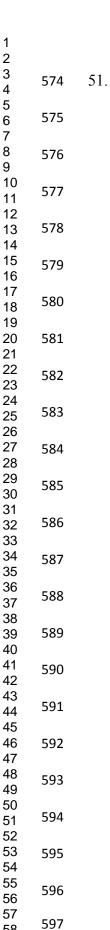
481          **13,** 788–794 (2014).

482    10.    Joly, P. *et al.* Prevalence projections of chronic diseases and impact of public health

483        intervention. *Biometrics* **69,** 109–17 (2013).

484    11.    Lee, Y. The recent decline in prevalence of dementia in developed countries: implications

485        for prevention in the Republic of Korea. *J. Korean Med. Sci.* **29,** 913–8 (2014).

486    12.    Barnes, D. E. & Yaffe, K. The projected effect of risk factor reduction on Alzheimer's

487        disease prevalence. *Lancet. Neurol.* **10,** 819–28 (2011).

488    13.    Manuel, D. G. *et al.* Projections of preventable risks for cardiovascular disease in Canada

489        to 2021: a microsimulation modelling approach. *C. Open* **2,** E94–E101 (2014).

490    14.    Tang, E. Y. H. *et al.* Current developments in dementia risk prediction modelling: An

491        updated systematic review. *PLoS One* **10,** 1–31 (2015).

492    15.    Stephan, B. C. M., Kurth, T., Matthews, F. E., Brayne, C. & Dufouil, C. Dementia risk

493        prediction in the population: are screening models accurate? *Nat. Rev. Neurol.* **6,** 318–326

494        (2010).

495    16.    Anstey, K. J. *et al.* A self-report risk index to predict occurrence of dementia in three

496        independent cohorts of older adults: The ANU-ADRI. *PLoS One* **9,** (2014).

497    17.    Wolfsgruber, S. *et al.* The CERAD neuropsychological assessment battery total score

498        detects and predicts alzheimer disease dementia with high diagnostic accuracy. *Am. J.*

499        *Geriatr. Psychiatry* **22,** 1017–1028 (2014).

500    18.    Barnes, D. E. *et al.* Predicting risk of dementia in older adults: The late-life dementia risk

501        index. *Neurology* **73,** 173–179 (2009).

502    19.    Chary, E. *et al.* Short- versus long-term prediction of dementia among subjects with low

503        and high educational levels. *Alzheimers. Dement.* **9,** 562–71 (2013).

504    20.    Jessen, F. *et al.* Prediction of Dementia in Primary Care Patients. *PLoS One* **6,** e16852

505   (2011).

506   21.   Song, X., Mitnitski, A. & Rockwood, K. Nontraditional risk factors combine to predict

507         Alzheimer disease and dementia. *Neurology* **77,** 227–234 (2011).

508   22.   Tierney, M. C., Moineddin, R. & McDowell, I. Prediction of all-cause dementia using

509         neuropsychological tests within 10 and 5 years of diagnosis in a community-based sample.

510         *J. Alzheimer's Dis.* **22,** 1231–1240 (2010).

511   23.   Meng, X., D'Arcy, C., Morgan, D. & Mousseau, D. D. Predicting the risk of dementia

512         among Canadian seniors: a useable practice-friendly diagnostic algorithm. *Alzheimer Dis.*

513         *Assoc. Disord.* **27,** 23–9 (2013).

514   24.   Kivipelto, M. *et al.* Risk score for the prediction of dementia risk in 20 years among

515         middle aged people: a longitudinal, population-based study. *Lancet Neurol.* **5,** 735–741

516         (2006).

517   25.   Reitz, C. *et al.* A Summary Risk Score for the Prediction of Alzheimer Disease in Elderly

518         Persons. *Arch. Neurol.* **67,** 835–41 (2010).

519   26.   Walters, K. *et al.* Predicting dementia risk in primary care: development and validation of

520         the Dementia Risk Score using routinely collected data. *BMC Med.* **14,** 1–12 (2016).

521   27.   Béland, Y. Canadian community health survey--methodological overview. *Heal. reports /*

522         *Stat. Canada, Can. Cent. Heal. Inf. = Rapp. sur la sant?? / Stat. Canada, Cent. Can.*

523         *d'information sur la sant??* **13,** 9–14 (2002).

524   28.   Jaakkimainen, R. L. *et al.* Identification of Physician-Diagnosed Alzheimer's Disease and

525         Related Dementias in Population-Based Administrative Data: A Validation Study Using

526         Family Physicians' Electronic Medical Records. *J. Alzheimer's Dis.* **54,** 337–349 (2016).

527   29.   Connolly, A., Gaehl, E., Martin, H., Morris, J. & Purandare, N. Underdiagnosis of

528    dementia in primary care: Variations in the observed prevalence and comparisons to the

529    expected prevalence. *Aging Ment. Health* **15,** 978–984 (2011).

530    30.   Kosteniuk, J. G. *et al.* Incidence and prevalence of dementia in linked administrative

531    health data in Saskatchewan, Canada: a retrospective cohort study. *BMC Geriatr.* **15,** 73

532    (2015).

533    31.   Harrell, F. E. *Regression Modeling Strategies with applicaitons to linear models, logistic

534    regression and survival analysis*. (Springer, 2001).

535    32.   Vergouwe, Y., Steyerberg, E. W., Eijkemans, M. J. C. & Habbema, J. D. F. Substantial

536    effective sample sizes were required for external validation studies of predictive logistic

537    regression models. *J. Clin. Epidemiol.* **58,** 475–483 (2005).

538    33.   Steyerberg, E. W. *Clinical Prediction Models*. (Springer, 2009).

539    34.   Harrell, F. E. *Hmisc: Harrell Miscellaneous. R package version 4.0-2*. (2016).

540    35.   R Core Team. R: A language and environment for statistical computing. (2016).

541    36.   Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. M. Transparent reporting of

542    a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The

543    TRIPOD Statement. *Eur. Urol.* **67,** 1142–1151 (2015).

544    37.   Exalto, L. G. *et al.* Midlife risk score for the prediction of dementia four decades later.

545    *Alzheimer's Dement.* **10,** 562–570 (2014).

546    38.   Exalto, L. G. *et al.* Risk score for prediction of 10 year dementia risk in individuals with

547    type 2 diabetes: A cohort study. *Lancet Diabetes Endocrinol.* **1,** 183–190 (2013).

548    39.   Taljaard, M. *et al.* Cardiovascular Disease Population Risk Tool (CVDPoRT): predictive

549    algorithm for assessing CVD risk in the community setting. A study protocol. *BMJ Open*

550    **4,** e006701 (2014).

551    40.    Manuel, D. G. *et al.* Measuring Burden of Unhealthy Behaviours Using a Multivariable

552           Predictive Approach: Life Expectancy Lost in Canada Attributable to Smoking, Alcohol,

553           Physical Inactivity, and Diet. *PLoS Med.* **13,** 1–27 (2016).

554    41.    Pampalon, R., Hamel, D., Gamache, P. & Raymond, G. A deprivation index for health

555           planning in Canada. *Chronic Dis. Can.* **29,** 178–91 (2009).

556    42.    Rubin, D. B. & Schenker, N. Multiple imputation in health-care databases: an overview

557           and some applications. *Stat. Med.* **10,** 585–98 (1991).

558    43.    Therneau, T. M., Grambsch, P. M. & Fleming, T. R. Martingale-based residuals for

559           survival models. *Biometrika* **77,** 147–160 (1990).

560    44.    Allison, P. D. *Survival analysis using SAS : A practical guide.* (SAS Institute, 2010).

561    45.    Van Houwelingen, J. C. & Le, C. S. Predictive value of statistical models. *Stat Med* **9,**

562           1303–1325 (1990).

563    46.    Ambler, G., Brady, A. R. & Royston, P. Simplifying a prognostic model: A simulation

564           study based on clinical data. *Stat. Med.* **21,** 3803–3822 (2002).

565    47.    Nagelkerke, N. J. D. A Note on a General Definition of the Coefficient of Determination.

566           *Biometrika* **78,** 691–692 (1991).

567    48.    Arkes, H. R. *et al.* The covariance decomposition of the probability score and its use in

568           evaluating prognostic estimates. SUPPORT Investigators. *Med. Decis. Making* **15,** 120–31

569           (1995).

570    49.    Cook, N. R. Statistical evaluation of prognostic versus diagnostic models: beyond the

571           ROC curve. *Clin. Chem.* **54,** 17–23 (2008).

572    50.    Cook, N. R. Comment: Measures to summarize and compare the predictive capacity of

573           markers. *Int. J. Biostat.* **6,** Article 22; discussion Article 25 (2010).

574    51.    Krause, J., Perer, A. & Bertini, E. Using Visual Analytics to Interpret Predictive Machine

575            Learning Models. *arXiv Prepr. arXiv1606.05685.* (2016).

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598 Table 1. Pre-specification of predictor variables for DemPoRT with initial degrees of freedom (df) allocation

| Variable | Scale | Initial Variable Specification | df |
|---|---|---|---|
| **Socio-demographic Factors** | | | |
| Age | Continuous | **5 knot spline:** Valid range: 55-102 (male), 55-101 (female) | 4 |
| Sex | Categorical | **Stratified:** Male; Female | NA |
| Ethnicity | Categorical | **7 categories:** Caucasian; African-American; Chinese; Aboriginal; Japanese/Korean/South East Asian/Filipino; Other/Multiple origin/Unknown/Latin American; South Asian/Arab/West Asian | 6 |
| Immigrant | Dichotomous | Yes; No | 1 |
| Education | Categorical | **4 categories:** Less than secondary school; Secondary school graduation; Some postsecondary; Postsecondary graduation | 3 |
| Marital Status | Categorical | **4 categories:** Now married/Common-law; Separated/Divorced; Widowed; Single | 3 |
| Neighborhood Social and Material Deprivation[41] | Ordinal | **3 categories:** Low (1st or 2nd quintile); High 4th or 5th quintile; Moderate (3rd quintile) | 2 |
| **General Health** | | | |
| Sense of belonging to local community | Ordinal | **4 categories:** Very strong; Somewhat strong; Somewhat weak; Very weak | 3 |
| Self-perceived stress | Ordinal | **5 categories:** Not at all stressful; Not very stressful; A bit stressful; Quite a bit stressful; Extremely stressful | 4 |
| Self-rated health | Ordinal | **5 categories:** Poor; Fair; Good; Very Good; Excellent | 4 |
| **Health Behaviors** | | | |
| Pack years of smoking | Continuous | **3 knot spline:** Valid range: 0-112 (male), 0-78 (female) | 2 |
| Smoking status | Categorical | **4 categories:** Non-smoker; Current smoker; Former smoker quit <5 years ago; Former smoker quit >5 years ago | 3 |
| Alcohol consumption (number of drinks last week) | Continuous | **3 knot spline:** Valid range: 0-50 (male), 0-24 (female) | 2 |
| Former drinker | Dichotomous | Yes; No | 1 |
| Consumption of fruit, salad, carrot and other vegetables (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-48 (male), 0-31 (female) | 2 |
| Potato consumption (average daily frequency) | Continuous | **3 knot spline:** Valid range: 0-2 | 2 |
| Juice consumption (average daily consumption | Continuous | **3 knot spline:** Valid range: 0-6 (male), 0-5 (female) | 2 |
| Leisure physical activity (average daily METs (kcal/kg/day)) | Continuous | **3 knot spline:** Valid range: 0-16 (male), 0-12 (female) | 2 |
| **Functional Measures** | | | |
| Personal hygiene and care | Dichotomous | Does not need help; Needs help | 1 |
| Locomotion in the home | Dichotomous | Does not need help; Needs help | 1 |
| Meal preparation | Dichotomous | Does not need help; Needs help | 1 |
| Running errands | Dichotomous | Does not need help; Needs help | 1 |
| Ordinary housework | Dichotomous | Does not need help; Needs help | 1 |
| Heavy housework | Dichotomous | Does not need help; Needs help | 1 |
| Finances | Dichotomous | Does not need help; Needs help | 1 |
| **Health Conditions** | | | |
| Heart disease | Dichotomous | Yes; No | 1 |
| Stroke | Dichotomous | Yes; No | 1 |
| Diabetes | Dichotomous | Yes; No | 1 |
| Mood disorder | Dichotomous | Yes; No | 1 |
| High blood pressure | Dichotomous | Yes; No | 1 |
| Body mass index | Continuous | **3 knot spline:** Valid range: 10-44 (male), 10-47 (female) | 2 |
| **Design** | | | |
| Survey year | Ordinal | **4 categories:** 2000/01, 2002/03, 2004/05, 2006/07 | 3 |

599 DemPoRT, Dementia Population Risk Tool; df, degrees of freedom; MET, metabolic equivalent task

# TRIPOD Checklist: Prediction Model Development and Validation

| Section/Topic | Item | | Checklist Item | Page |
|---|---|---|---|---|
| **Title and abstract** | | | | |
| Title | 1 | D;V | Identify the study as developing and/or validating a multivariable prediction model, the target population, and the outcome to be predicted. | 3 |
| Abstract | 2 | D;V | Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions. | 3 |
| **Introduction** | | | | |
| Background and objectives | 3a | D;V | Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models. | 5-7 |
| | 3b | D;V | Specify the objectives, including whether the study describes the development or validation of the model or both. | 7,8 |
| **Methods** | | | | |
| Source of data | 4a | D;V | Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable. | 8,9 |
| | 4b | D;V | Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up. | 8-10 |
| Participants | 5a | D;V | Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres. | 8,9 |
| | 5b | D;V | Describe eligibility criteria for participants. | 8-10 |
| | 5c | D;V | Give details of treatments received, if relevant. | NA |
| Outcome | 6a | D;V | Clearly define the outcome that is predicted by the prediction model, including how and when assessed. | 9,10 |
| | 6b | D;V | Report any actions to blind assessment of the outcome to be predicted. | 12 |
| Predictors | 7a | D;V | Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured. | 11,12, 25 |
| | 7b | D;V | Report any actions to blind assessment of predictors for the outcome and other predictors. | 12 |
| Sample size | 8 | D;V | Explain how the study size was arrived at. | 8-12 |
| Missing data | 9 | D;V | Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method. | 12,13 |
| Statistical analysis methods | 10a | D | Describe how predictors were handled in the analyses. | 13,14, 25 |
| | 10b | D | Specify type of model, all model-building procedures (including any predictor selection), and method for internal validation. | 14,15 |
| | 10c | V | For validation, describe how the predictions were calculated. | 16 |
| | 10d | D;V | Specify all measures used to assess model performance and, if relevant, to compare multiple models. | 16 |
| | 10e | V | Describe any model updating (e.g., recalibration) arising from the validation, if done. | NA |
| Risk groups | 11 | D;V | Provide details on how risk groups were created, if done. | 16 |
| Development vs. validation | 12 | V | For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors. | 16,17 |
| **Results** | | | | |
| Participants | 13a | D;V | Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful. | 10 |
| | 13b | D;V | Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome. | NA |
| | 13c | V | For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome). | NA |
| Model development | 14a | D | Specify the number of participants and outcome events in each analysis. | NA |
| | 14b | D | If done, report the unadjusted association between each candidate predictor and outcome. | NA |
| Model specification | 15a | D | Present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point). | NA |
| | 15b | D | Explain how to the use the prediction model. | NA |
| Model performance | 16 | D;V | Report performance measures (with CIs) for the prediction model. | NA |
| Model-updating | 17 | V | If done, report the results from any model updating (i.e., model specification, model performance). | NA |
| **Discussion** | | | | |
| Limitations | 18 | D;V | Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data). | 4,17 |
| Interpretation | 19a | V | For validation, discuss the results with reference to performance in the development data, and any other validation data. | NA |
| | 19b | D;V | Give an overall interpretation of the results, considering objectives, limitations, results from similar studies, and other relevant evidence. | NA |
| Implications | 20 | D;V | Discuss the potential clinical use of the model and implications for future research. | 5-7 |
| **Other information** | | | | |
| Supplementary information | 21 | D;V | Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets. | 4,18/ 19 |
| Funding | 22 | D;V | Give the source of funding and the role of the funders for the present study. | 20 |

\*Items relevant only to the development of a prediction model are denoted by D, items relating solely to a validation of a prediction model are denoted by V, and items relating to both are denoted D;V. We recommend using the TRIPOD Checklist in conjunction with the TRIPOD Explanation and Elaboration document.