

BMJ Open Dementia Population Risk Tool (DemPoRT): study protocol for a predictive algorithm assessing dementia risk in the community

Stacey Fisher,^{1,2,3} Amy Hsu,^{1,2} Nassim Mojaverian,² Monica Taljaard,^{1,3} Gregory Huyer,⁴ Douglas G Manuel,^{1,2,3,5,6} Peter Tanuseputro^{1,2,6,7,8}

To cite: Fisher S, Hsu A, Mojaverian N, *et al.* Dementia Population Risk Tool (DemPoRT): study protocol for a predictive algorithm assessing dementia risk in the community. *BMJ Open* 2017;7:e018018. doi:10.1136/bmjopen-2017-018018

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2017-018018>).

Received 31 May 2017
Revised 25 August 2017
Accepted 14 September 2017



CrossMark

For numbered affiliations see end of article.

Correspondence to
Dr Peter Tanuseputro;
ptanuseputro@ohri.ca

ABSTRACT

Introduction The burden of disease from dementia is a growing global concern as incidence increases dramatically with age, and average life expectancy has been increasing around the world. Planning for an ageing population requires reliable projections of dementia prevalence; 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 modelling techniques and will be used to project dementia prevalence.

Methods and analysis The derivation cohort will consist of elderly Ontario respondents of the Canadian Community Health Survey (CCHS) (2001, 2003, 2005 and 2007; 18 764 males and 25 288 females). Prespecified predictors include sociodemographic, general health, behavioural, functional and health condition variables. Incident dementia will be identified through individual linkage of survey respondents to population-level administrative healthcare databases (1797 and 3281 events, and 117 795 and 166 573 person-years of follow-up, for males and females, respectively, until 31 March 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 4600 males and 6300 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 Research ethics approval has been granted 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, pre-results.

INTRODUCTION

The burden of disease from dementia is a growing global concern as incidence increases dramatically with age, and average

Strengths and limitations of this study

- The Dementia Population Risk Tool (DemPoRT) will be developed and validated using predictors from large population-based community health surveys that are individually linked to routinely collected health administration data in Ontario. To our knowledge, DemPoRT will be the first algorithm designed to predict and project dementia incidence at the population level.
- Although repeat predictor assessment and detailed cognitive testing to ascertain dementia diagnoses is preferable, it is not available or feasible at the population level.
- Statistical overfitting is a concern; however, full prespecification of the analysis plan and predictors will limit this risk.
- Although a rigorous approach to model development will be used, further validation will be needed to assess generalisability, and calibration will be required for application in other jurisdictions.
- DemPoRT will be used to produce improved estimates of future dementia burden, will assess the contribution of specific risk factors to the population risk and will identify population subgroups at high risk of developing dementia. This information will be used by policy makers to prepare for and reduce dementia impact.

life expectancy has been increasing around the world.^{1 2} Planning for an ageing population requires reliable projections of dementia burden and the implications for resource requirements. Existing population-level projections for dementia, however, are overly simplistic and likely inaccurate.³

Limitations of current dementia projection methodology

Almost all existing dementia projections have used extrapolation and macrosimulation methods, which are simplistic and make assumptions that may not hold true into the future.³ Most extrapolations simply

apply current age-specific and sex-specific prevalence estimates of dementia to future population projections. Macrosimulations typically use estimates of dementia incidence and mortality, stratified by age and sex, to simulate disease prevalence as the population ages.¹⁻⁴⁻⁶ Projections from extrapolations incorrectly assume that the risk of mortality among those with and without dementia are equivalent,^{7,8} and both methods assume that the age-specific and sex-specific prevalence of dementia risk factors will not change with time. The assumption of stable risk factor prevalence is widely thought to be the major source of error in existing dementia projections.^{3,9-11}

Changing trends of dementia risk factors has the potential to have a dramatic impact on dementia prevalence estimates, as up to 50% of dementia cases have been attributed to modifiable factors,^{9,12} and the prevalence of several factors has been projected to change significantly in the near future. For example, the population prevalence of diabetes and obesity in Canada has been projected to increase, while smoking, hypertension and dyslipidaemia have been projected to decline.¹³ Consideration of risk factor prevalence is therefore important to improve the accuracy of dementia projections, and simple extrapolations and macrosimulations are often inadequate.

Predictive multivariable modelling of dementia incidence

Population-based predictive risk algorithms examine the effect of risk factors on dementia incidence and can be used for dementia burden projection. Population-based data that contain detailed risk factor information, such as health surveys, are linked at the individual level to administrative data that capture dementia development. A multivariable model of dementia incidence is derived, validated against external data, and predictive performance is assessed. Once developed, the algorithm can be used to project disease incidence and prevalence. To obtain prevalence projections, the algorithm can be integrated in to a microsimulation model such as Statistics Canada's Population Health Models (POHEM). POHEM dynamically models individual life trajectories of a population representative of Canada including births, deaths and migration, disease incidence and progression and exposure to risk factors, facilitating detailed examination of the influence of changing risk factor prevalence on future dementia prevalence.

Predictive risk algorithms can also be used to describe the risk of dementia in the population, assess the contribution of specific risk factors to the population risk, identify high-risk groups and evaluate risk reduction strategies.

Existing dementia prediction models

Many models have been developed to predict risk of dementia,¹⁴⁻²⁶ most with the primary goal of identifying individuals in the clinical setting at high risk. They have varying discriminative ability (c-statistics ranging from 0.49¹⁶ to 0.89¹⁷) and have generally been derived from small samples, rarely including more than a few

thousand individuals. Existing models are therefore simplistic, including few predictors and rarely including interaction or non-linear terms, facilitating understanding and use by physicians in clinical practice but limiting discriminatory ability and predictive accuracy. Walters *et al*²⁶ developed an algorithm for predicting 5-year dementia risk among individuals 60-79 years of age in the UK using an enormous derivation dataset of 800 000 individuals and a simple model. The derivation model had a c-statistic of 0.84 (95% CI 0.81 to 0.87), but a low positive predictive value at most risk thresholds, and therefore is poor at identifying those at high risk of dementia. Additionally, as most dementia risk models are intended for use in the clinical setting, many include results from neuropsychological tests,¹⁷⁻²³ MRI findings¹⁸ and apolipoprotein E (APOE) genotype.^{18,24,25} The inclusion of these variables, however, limits the application of these models as these variables are not available at the population level.

The objective of this study is to develop and validate the Dementia Population Risk Tool (DemPoRT) algorithm to predict dementia incidence in the population setting. This will be done using multivariable modelling techniques, linking self-reported risk factors captured by a population-based health survey in Canada with administrative databases across healthcare sectors that capture healthcare-diagnosed dementia. DemPoRT will be developed using a large population-based dataset using only variables that are available at the population level, allowing for population-level application. DemPoRT will also use many methodological improvements over existing models. This protocol prespecifies the predictor variables and analytic plan for model development, reducing the potential for overfitting and bias and improving transparency. Interaction terms and flexible functions for continuous predictors will be investigated, increasing potential discriminative ability. The prespecified analytic plan avoids data-driven variable selection procedures, further reducing the potential for bias.

To our knowledge, the DemPoRT predictive model will be the first algorithm designed to predict and project dementia incidence at the population level. It will be used to estimate the future burden of dementia using techniques that consider changes in risk factor prevalence and will identify modifiable risk factors that can be targeted by individuals, clinicians and policy makers to reduce the burden of dementia.

METHODS AND ANALYSIS

Study design

Two DemPoRT models, one for males and females, will be derived and validated using population-based data in Ontario, Canada, a multicultural province with 13.6 million residents. Predictors will be obtained from the Canadian Community Health Surveys (CCHS), and outcomes (ie, diagnosis of dementia) will be obtained from routinely collected healthcare data.

The derivation cohorts will consist of eligible respondents of the 2001, 2003, 2005 and 2007 CCHS (cycles 1.1, 2.1, 3.1 and 4.1), while validation cohorts will consist of respondents to the 2008/2009 cycle. The CCHS is a national, cross-sectional survey developed by Statistics Canada to collect information related to health and healthcare utilisation of the Canadian population. The survey has a multistage, stratified cluster design that represents approximately 98% of the Canadian population aged 12 years and over and attained an average response rate of 79% over the study period. The CCHS is conducted through telephone and in-person interviews, and all responses are self-reported. The details of survey methodology have been published elsewhere.²⁷ Survey respondents will be excluded if they are less than 55 years of age at survey administration, self-reported a history of dementia or are not eligible for Ontario's universal health insurance. If a respondent was included in more than one CCHS cycle, only their earliest survey response will be used.

Outcome

Survey respondents diagnosed with dementia will be identified through individual linkage to several population-based administrative databases at the Institute for Clinical Evaluative Sciences. Dementia case ascertainment is based on a validated definition: one hospital record OR three physician claim records at least 30 days apart within a 2-year period OR a dispensing record for a cholinesterase inhibitor from Ontario Drug Benefit. This definition has a sensitivity of 79.3% and a specificity of 99.1% when validated against emergency medical record data.²⁸ Due to known underdiagnosis of dementia,^{29,30} we will supplement this definition by adding survey respondents with dementia codes captured on home care and long-term care assessments (dementia flag AND Cognitive Performance Scale score ≥ 2) using the Resident Assessment Instrument-Home Care database and the Continuing Care Reporting System, respectively. We have found this addition adds substantially (approximately 18%) to the number of dementia cases captured.

Survey respondents with dementia will be excluded if they meet the criteria for dementia within 2 years of survey administration (to remove potentially prevalent cases) or are younger than 65 years of age at the time of dementia diagnosis (to exclude early onset dementia which likely has a different set of risk factors). Eligible survey respondents will be followed from the date of survey administration or age 65 years, whichever came later, until the earliest date of: dementia ascertainment, death (defined as competing risk), loss to follow-up (defined as loss of healthcare eligibility) or end of study (31 March 2014).

Sample size

The male and female derivation cohorts consist of 18 764 and 25 288 respondents, and 117 795 and 166 573 person-years of follow-up, respectively. For predictive models with time to event outcomes, the number of participants

experiencing the event should exceed 10 times the number of degrees of freedom (df) to ensure adequate sample size.³¹ The number of dementia events in the derivation cohort is 1797 for men and 3281 for women; therefore, the maximum number of total df for each of the DemPoRT models is 179 and 328, respectively, which we do not anticipate surpassing.

The validation cohorts will consist of approximately 4600 males and 6300 females, and 15 000 and 21 000 person-years of follow-up, respectively. Vergouwe *et al*³² recommend a minimum of 100 events and 100 non-events for external validation studies. We expect approximately 225 events for men and 400 for women in our validation cohort.

Analysis plan

The analysis plan was developed following guidelines by Harrell³¹ and Steyerberg³³ after accessing the derivation data set but prior to model fitting or descriptive analyses involving exposure-outcome associations. This was done to avoid type 1 error introduced by data-driven variable selection or model specification. Key considerations of our analysis approach include full prespecification of the predictor variables, use of flexible functions for continuous predictors and preserving statistical properties by avoiding data-driven variable selection procedures. Analysis will be conducted using Harrell's Hmisc³⁴ package of functions in R³⁵ as well as SAS V.9.4.

This study protocol and the reporting of our model estimation results will be guided by the TRIPOD statement for multivariable predictive models.³⁶

Identification of predictors

Predictor variables were identified through review of existing predictive algorithms for dementia^{9,16,18-22,24-26,37,38} and comparison with available data collected in the CCHS. Variable inclusion was informed by consultation with subject-matter experts and the project's advisory team and informed by our previous work developing predictive models for cardiovascular disease and life expectancy.^{39,40}

Variables with narrow distributions or insufficient variation were excluded. Obvious cases of redundancy (eg, alternate definitions of the same underlying behaviour) were not included. A total of 32 predictor variables were identified: seven sociodemographic, three general health, nine behavioural, seven functional, five health conditions and one design variable (CCHS survey cycle). As the effect of dementia risk factors varies by sex, separate models will be derived for men and women. Education, rather than individual income, was selected as a predictor due to several concerns with income including lack of generalisability, measurement error, stability over time and substantial missing values. Neighbourhood social and material deprivation is captured using Pampalon's deprivation index.⁴¹ Indicator variables for smoking status were created to allow the inclusion of smoking pack-years as a continuous predictor. The models will additionally include age interactions with the behavioural, functional

and health condition variables as the effect of these risk factors on dementia are expected to vary with age. Detailed definitions and measurement of the predictor variables are presented in [table 1](#).

Data cleaning and coding of predictors

Continuous variables will be inspected using boxplots and descriptive statistics to determine values outside a plausible range. Values that are clearly erroneous will be corrected, where possible, or set to missing. Continuous predictors with highly skewed distributions will be truncated to the 99.5th percentile. Categorisation of continuous variables will be avoided to minimise the loss of predictive information. All data cleaning and coding will occur prior to examining exposure-outcome associations.

Missing data

As traditional complete cases analyses suffer from inefficiency, selection bias and other limitations,³³ multiple imputation methods will be used to impute missing values using the 'aregImpute' function in the HMisc library.³⁴ This function simultaneously imputes missing values using predictive mean matching and uses bootstrapping to take all aspects of uncertainty into account. The imputation model will consist of the full list of predictor variables, time to event and censoring variables, as well as auxiliary variables that are not predictors but may nevertheless be useful in generating imputed values (eg, income). The final model will be estimated in each of five multiple imputation data sets, and the results will be combined using the rules developed by Rubin and Schenker⁴² to account for imputation uncertainty.

Model specification

Initial sex-specific main effects models will be fit using the prespecified predictors and an initial degree of freedom allocation for each predictor ([table 1](#)). Decisions on initial degree of freedom allocations were informed by the anticipated importance of each predictor and known dose-response relationships with dementia. Continuous predictors will be flexibly modelled using restricted cubic splines, with the knots placed at fixed quantiles of the distribution (eg, 5th, 27.5th, 50th, 72.5th and 95th centiles). Frequency distributions for categorical predictors will be examined, and categories with small numbers of respondents will be combined, with analysts blinded to the number of events per category, to avoid instability in the regression analyses. Ordinal variables will be specified as either linear terms or as categorical if the expected association is more complex. Interactions will be restricted to linear terms. The initial model specification, presented in [table 1](#), includes a total of 86 df (63 main, 23 interaction).

Partial association χ^2 statistics for each predictor minus their df (to level the playing field among predictors with varying df) will be plotted in descending order. Variables with higher predictive potential will retain their initial df, while predictors with lower predictive potential will be modelled as simple linear terms or recoded by combining

infrequent categories. This process of model specification does not increase the type 1 error rate, because all predictors will be retained in the full model regardless of their strength of association.³¹

Model estimation

The initial models will be estimated using competing risk Cox proportional hazards regression with time to dementia ascertainment as the outcome and death as a competing risk. Alternative model specifications, including subdistribution hazard and flexible parametric models, will be considered. All predictors will be centred about their means. A formal check of multicollinearity will be carried out using a variable clustering algorithm.³¹

Proportional hazards models assume that the relative risk of the outcome between strata of predictors and the baseline risk must be constant over time. Violation of this assumption has been shown to produce biased results,⁴³ although it has also been argued that the estimated coefficients of time-varying variables can simply be interpreted as an average rather than instantaneous hazard.⁴⁴ Plots of raw and smoothed scaled Schoenfeld residuals versus time for each predictor will be assessed to test this assumption and identify non-proportionality. If a violation of this assumption is identified, we will consider addition of interaction terms between the predictor and log-transformed time.

Although the risk of overfitting will be minimal due to prespecification of the models and a large sample size, the need for overfitting adjustment will be assessed. The degree of overfitting will be estimated using the heuristic shrinkage estimator, based on the log likelihood ratio χ^2 statistic for the full model.⁴⁵ If shrinkage is <0.90 , models will be adjusted for overfitting.

Estimation of the reduced models

Model prespecification has advantages in limiting overfitting and spurious statistical significance but can result in a final model that is overly complex, difficult to interpret and difficult to apply. Unnecessary predictor variables also distort the estimated effects of other predictors making the model more computationally intensive. It is suggested that a more parsimonious model that retains most of the prognostic information and performs as well as or better than the full model can be derived without increasing the type 1 error rate.^{31 46} We will identify a more parsimonious model using a stepdown procedure described by Ambler *et al*,⁴⁶ which involves deleting the variable that results in the smallest decrease in model R^2 until removal leads to an R^2 that is less than a desired level. The reduced model will be evaluated against the full model using Akaike's Information Criterion and by examining the effect on discrimination and calibration.

DemPoRT will be developed and validated using temporal split samples; however, the final regression coefficients will use the full data set to maximise follow-up duration. A cohort-specific intercept and/or interaction term may be included in the final model if the

Table 1 Prespecification of predictor variables for DemPoRT with initial degrees of freedom (df) allocation

Variable	Scale	Initial variable specification	df
Sociodemographic factors			
Age	Continuous	5 knot spline: valid range: 55–102 (male), 55–101 (female)	4
Sex	Categorical	Stratified: Male; female	NA
Ethnicity	Categorical	Seven 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	Four categories: less than secondary school; secondary school graduation; some postsecondary; postsecondary graduation	3
Marital status	Categorical	Four categories: now married/common law; separated/divorced; widowed; single	3
Neighbourhood social and material Deprivation ⁴¹	Ordinal	Three categories: low (1st or 2nd quintile); high 4th or 5th quintile; moderate (3rd quintile)	2
General health			
Sense of belonging to local community	Ordinal	Four categories: very strong; somewhat strong; somewhat weak; very weak	3
Self-perceived stress	Ordinal	Five categories: not at all stressful; not very stressful; a bit stressful; quite a bit stressful; extremely stressful	4
Self-rated health	Ordinal	Five categories: poor; fair; good; very good; excellent	4
Health behaviours			
Pack years of smoking	Continuous	3 knot spline: valid range: 0–112 (male), 0–78 (female)	2
Smoking status	Categorical	Four 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

Continued

Table 1 Continued

Variable	Scale	Initial variable specification	df
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	Four categories: 2000/2001, 2002/2003, 2004/2005, 2006/2007	3

DemPoRT, Dementia Population Risk Tool; METs, metabolic equivalent tasks.

derivation and validation cohorts differ; otherwise, the final combined model will maintain the same predictors and form as the derivation model.

Assessment of predictive performance

Predictive performance in the derivation and validation cohorts will be assessed and reported using overall measures of predictive accuracy, discrimination and calibration. Accuracy will be assessed with Nagelkerke's R^2 ⁴⁷ and the Brier score.⁴⁸ Discrimination will be assessed using the concordance statistic. Model calibration is especially important in the development of prognostic models, as probabilities of future risk are of primary interest.^{33 49 50} Calibration will be assessed by comparing the observed and predicted risk of dementia within vigintiles (20 groups of equal frequency) of predicted risk with emphasis on visual inspection of plots rather than formal statistical significance testing, which can be influenced by large sample sizes.³² Calibration slopes will be generated by regressing the outcome in the validation cohort on the predicted dementia risk, reflecting the combined effect of overfitting to the derivation data as well as true differences in effects of predictors. Deviation of the slope from 1 (perfect calibration) will be tested using a Wald or likelihood ratio test. Calibration within predefined subgroups of importance to clinicians and policy makers (eg, age group, health behaviour, sociodemographic groups and health conditions) will additionally be evaluated. The clinically relevant standard of calibration was defined as less than 20% difference between observed and predicted estimates within subgroups with a dementia prevalence of at least 5%. All model performance measures will be calculated using the first of the multiply imputed data sets.

Model presentation

The final regression model, derived from the combined sample of the derivation and validation cohorts, will be presented using estimated HRs and 95% CIs, along with results for the derivation and validation cohorts separately. We have found, however, this usual presentation less meaningful when presenting complex models.³⁹ To allow interpretation of the estimated effect of each predictor, model behaviour will additionally be described using interactive visual tools to display the shape of the effect of each predictor.⁵¹ The regression formula will

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

Analyses beyond initial model development

We will conduct further analyses exploring the added predictive ability of novel risk factors that were ascertained in single CCHS cycles (eg, sedentary activity, cognitive stimulation, sleep quality and duration and deafness), as well as risk factors that can be ascertained through linkage of additional data sources and similar cohorts (eg, air pollution, detailed dietary consumption, lipid levels and blood pressure). In addition, sensitivity analysis of the age at survey administration cut-off used for cohort creation will be performed.

Once developed, DemPoRT will be used to project dementia incidence under different assumptions by entering counterfactual risk factor levels into the algorithm at the population level, or at individual level and summed, and will be integrated in to POHEM for micro-simulation modelling of prevalence projections.

A second, causal model (DemPoRT-C) will also be created to assess the relative contribution of lifestyle, sociodemographic and health factors to dementia incidence. Development will exclude variables believed to be in the causal pathway of dementia occurrence (eg, self-rated health and functional measures) to reduce the attenuation of hazards from upstream risk factors but will otherwise be the same as in the predictive model. DemPoRT-C will be applied to the most recent 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 prespecification 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 generalisability, 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.

DemPoRT results will be submitted for publication in peer-review journals and presented at scientific meetings. A web-based individual-level calculator will be created if the models are appropriate for individual use. Although DemPoRT emphasises risk prediction at the population level, we have found that individual-level calculators are an effective engagement and translation tool for both the general public and knowledge users.

CONCLUSIONS

To the best of our knowledge, DemPoRT will be the first population-based algorithm designed to predict and project dementia incidence at the population level. The DemPoRT models will produce estimates of future dementia burden that we believe will be more accurate than existing estimates, will assess the contribution of specific risk factors to the population risk and will identify groups at high risk of developing dementia.

Author affiliations

¹Clinical Epidemiology Program, Ottawa Hospital Research Institute, Ottawa, Ontario, Canada

²Institute for Clinical Evaluative Sciences, Ottawa, Ontario, Canada

³Department of Epidemiology and Community Medicine, University of Ottawa, Ottawa, Ontario, Canada

⁴Telfer School of Management, University of Ottawa, Ottawa, Ontario, Canada

⁵Statistics Canada, Ottawa, Ontario, Canada

⁶Department of Family Medicine, University of Ottawa, Ottawa, Ontario, Canada

⁷Department of Medicine, University of Ottawa, Ottawa, Ontario, Canada

⁸Bruyère Research Institute, Ottawa, Ontario, Canada

Contributors SF drafted and revised the manuscript and contributed to the study design and protocol development. NM contributed to the study design, protocol development and provided data/statistical support. AH, MT, DM and GH contributed to the design of the study and protocol development. PT is the lead investigator of the study and was responsible for the conception of the project, the grant application study design and protocol development. All authors provided critical reviews of the manuscript and reviewed the final version.

Funding This work is supported by the Canadian Institutes of Health Research operating grant MOP142237. SF is supported by a doctoral award from the Canadian Institutes of Health Research. This project is supported by the Institute for Clinical Evaluative Sciences (ICES), which is funded by an annual grant from

the Ontario Ministry of Health and Long-Term Care (MOHLTC). The sponsors have no role in the design or conduct of the study; in the collection, analysis or interpretation of the data; or in the preparation, review or approval of the manuscript. The options, results and conclusions reported are those of the authors and are independent from the funding source. No endorsement by ICES or the Ontario MOHLTC is intended or should be inferred.

Competing interests None declared.

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

Provenance and peer review Not commissioned; externally peer reviewed.

Open Access This is an Open Access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>

© Article author(s) (or their employer(s) unless otherwise stated in the text of the article) 2017. All rights reserved. No commercial use is permitted unless otherwise expressly granted.

REFERENCES

- 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* 1998;88:1337–42.
- Brayne C. The elephant in the room - healthy brains in later life, epidemiology and public health. *Nat Rev Neurosci* 2007;8:233–9.
- Norton S, Matthews FE, Brayne C. A commentary on studies presenting projections of the future prevalence of dementia. *BMC Public Health* 2013;13:1.
- Sloane PD, Zimmerman S, Suchindran C, et al. The public health impact of Alzheimer's disease, 2000-2050: potential implication of treatment advances. *Annu Rev Public Health* 2002;23:213–31.
- Mura T, Dartigues JF, Berr C. How many dementia cases in France and Europe? Alternative projections and scenarios 2010-2050. *Eur J Neurol* 2010;17:252–9.
- Hebert LE, Scherr PA, Bienias JL, et al. Alzheimer disease in the US population: prevalence estimates using the 2000 census. *Arch Neurol* 2003;60:1119–22.
- Brookmeyer R, Johnson E, Ziegler-Graham K, et al. Forecasting the global burden of Alzheimer's disease. *Alzheimers Dement* 2007;3:186–91.
- Dewey ME, Chen CM. Neurosis and mortality in persons aged 65 and over living in the community: a systematic review of the literature. *Int J Geriatr Psychiatry* 2004;19:554–7.
- Norton S, Matthews FE, Barnes DE, et al. Potential for primary prevention of Alzheimer's disease: an analysis of population-based data. *Lancet Neurol* 2014;13:788–94.
- Joly P, Touraine C, Georget A, et al. Prevalence projections of chronic diseases and impact of public health intervention. *Biometrics* 2013;69:109–17.
- Lee Y. The recent decline in prevalence of dementia in developed countries: implications for prevention in the Republic of Korea. *J Korean Med Sci* 2014;29:913–8.
- Barnes DE, Yaffe K. The projected effect of risk factor reduction on Alzheimer's disease prevalence. *Lancet Neurol* 2011;10:819–28.
- Manuel DG, Tuna M, Hennessy D, et al. Projections of preventable risks for cardiovascular disease in Canada to 2021: a microsimulation modelling approach. *CMAJ Open* 2014;2:E94–101.
- Tang EY, Harrison SL, Errington L, et al. Current developments in dementia risk prediction modelling: an updated systematic review. *PLoS One* 2015;10:e0136181–31.
- Stephan BC, Kurth T, Matthews FE, et al. Dementia risk prediction in the population: are screening models accurate? *Nat Rev Neurol* 2010;6:318–26.
- Anstey KJ, Cherbuin N, Herath PM, et al. A self-report risk index to predict occurrence of dementia in three independent cohorts of older adults: the ANU-ADRI. *PLoS One* 2014;9:e86141.
- Wolfsgruber S, Jessen F, Wiese B, et al. The CERAD neuropsychological assessment battery total score detects and predicts Alzheimer disease dementia with high diagnostic accuracy. *Am J Geriatr Psychiatry* 2014;22:1017–28.
- Barnes DE, Covinsky KE, Whitmer RA, et al. Predicting risk of dementia in older adults: the late-life dementia risk index. *Neurology* 2009;73:173–9.

19. Chary E, Amieva H, Pérès K, *et al.* Short- versus long-term prediction of dementia among subjects with low and high educational levels. *Alzheimers Dement* 2013;9:562–71.
20. Jessen F, Wiese B, Bickel H, *et al.* AgeCoDe Study Group. Prediction of dementia in primary care patients. *PLoS One* 2011;6:e16852.
21. Song X, Mitnitski A, Rockwood K. Nontraditional risk factors combine to predict Alzheimer disease and dementia. *Neurology* 2011;77:227–34.
22. Tierney MC, Moineddin R, McDowell I. Prediction of all-cause dementia using neuropsychological tests within 10 and 5 years of diagnosis in a community-based sample. *J Alzheimers Dis* 2010;22:1231–40.
23. Meng X, D'Arcy C, Morgan D, *et al.* Predicting the risk of dementia among Canadian seniors: a useable practice-friendly diagnostic algorithm. *Alzheimer Dis Assoc Disord* 2013;27:23–9.
24. Kivipelto M, Ngandu T, Laatikainen T, *et al.* Risk score for the prediction of dementia risk in 20 years among middle aged people: a longitudinal, population-based study. *Lancet Neurol* 2006;5:735–41.
25. Reitz C, Tang MX, Schupf N, *et al.* A summary risk score for the prediction of Alzheimer disease in elderly persons. *Arch Neurol* 2010;67:835–41.
26. Walters K, Hardoon S, Petersen I, *et al.* Predicting dementia risk in primary care: development and validation of the Dementia Risk Score using routinely collected data. *BMC Med* 2016;14:1–12.
27. Béland Y. Canadian community health survey-methodological overview. Heal. reports / Statanada, Can. Cent. Heal. Inf. = Rapp. sur la sant?? / Stat. Canada, Cen.Ct. Can. d'information sur la sant?? *Health Rep* 2002;13:9–14.
28. Jaakkimainen RL, Bronskill SE, Tierney MC, *et al.* Identification of physician-diagnosed alzheimer's disease and related dementias in population-based administrative data: A validation study using family physicians' electronic medical records. *J Alzheimers Dis* 2016;54:337–49.
29. Connolly A, Gaehl E, Martin H, *et al.* Underdiagnosis of dementia in primary care: variations in the observed prevalence and comparisons to the expected prevalence. *Aging Ment Health* 2011;15:978–84.
30. Kosteniuk JG, Morgan DG, O'Connell ME, *et al.* Incidence and prevalence of dementia in linked administrative health data in Saskatchewan, Canada: a retrospective cohort study. *BMC Geriatr* 2015;15:73.
31. Harrell FE. Regression Modeling Strategies with applicaitons to linear models, logistic regression and survival analysis. *Springer* 2001.
32. Vergouwe Y, Steyerberg EW, Eijkemans MJ, *et al.* Substantial effective sample sizes were required for external validation studies of predictive logistic regression models. *J Clin Epidemiol* 2005;58:475–83.
33. Steyerberg EW. *Clinical Prediction Models*. New York: Springer, 2009.
34. Harrell FE. *Hmisc: harrell miscellaneous.R package version*, 2016. 4.0-2.
35. Core Team R. *R: A language and environment for statistical computing*, 2016.
36. Collins GS, Reitsma JB, Altman DG, *et al.* Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *Eur Urol* 2015;67:1142–51.
37. Exalto LG, Quesenberry CP, Barnes D, *et al.* Midlife risk score for the prediction of dementia four decades later. *Alzheimers Dement* 2014;10:562–70.
38. Exalto LG, Biessels GJ, Karter AJ, *et al.* Risk score for prediction of 10 year dementia risk in individuals with type 2 diabetes: a cohort study. *Lancet Diabetes Endocrinol* 2013;1:183–90.
39. Taljaard M, Tuna M, Bennett C, *et al.* Cardiovascular Disease Population Risk Tool (CVDPORT): predictive algorithm for assessing CVD risk in the community setting. A study protocol. *BMJ Open* 2014;4:e006701.
40. Manuel DG, Perez R, Sanmartin C, *et al.* Measuring burden of unhealthy behaviours using a multivariable predictive approach: life expectancy lost in Canada attributable to smoking, alcohol, physical inactivity, and diet. *PLoS Med* 2016;13:e1002082:27.
41. Pampalon R, Hamel D, Gamache P, *et al.* A deprivation index for health planning in Canada. *Chronic Dis Can* 2009;29:178–91.
42. Rubin DB, Schenker N. Multiple imputation in health-care databases: an overview and some applications. *Stat Med* 1991;10:585–98.
43. Therneau TM, Grambsch PM, Fleming TR. Martingale-based residuals for survival models. *Biometrika* 1990;77:147–60.
44. Allison PD. Survival analysis using SAS : a practical guide. *SAS Institute* 2010.
45. Van Houwelingen JC, Le Cessie S. Predictive value of statistical models. *Stat Med* 1990;9:1303–25.
46. Ambler G, Brady AR, Royston P. Simplifying a prognostic model: a simulation study based on clinical data. *Stat Med* 2002;21:3803–22.
47. Nagelkerke NJD. A note on a general definition of the coefficient of determination. *Biometrika* 1991;78:691–2.
48. Arkes HR, Dawson NV, Speroff T, *et al.* The covariance decomposition of the probability score and its use in evaluating prognostic estimates. SUPPORT Investigators. *Med Decis Making* 1995;15:120–31.
49. Cook NR. Statistical evaluation of prognostic versus diagnostic models: beyond the ROC curve. *Clin Chem* 2008;54:17–23.
50. Cook NR. Comment: Measures to summarize and compare the predictive capacity of markers. *Int J Biostat* 2010;6:Article 22; discussion Article 25.
51. Krause J, Perer A, Bertini E. Using Visual Analytics to Interpret Predictive Machine Learning Models *arXiv Prepr. arXiv1606* 2016;05685.