**Distribution Agreement**

In presenting this thesis as a partial fulfillment of the requirements for a degree from Emory University, I hereby grant to Emory University and its agents the non-exclusive license to archive, make accessible, and display my thesis in whole or in part in all forms of media, now or hereafter now, including display on the World Wide Web. I understand that I may select some access restrictions as part of the online submission of this thesis. I retain all ownership rights to the copyright of the thesis. I also retain the right to use in future works (such as articles or books) all or part of this thesis.

Xizhu (Emilia) Liu                                                                          April 2, 2022

Alzheimer's Disease Pathology Imputation and Risk Prediction Using Clinical Indices

by

Xizhu (Emilia) Liu

Jingjing Yang, Ph.D.
Adviser

Department of Quantitative Theory and Methods

Jingjing Yang, Ph.D.

Adviser

Zhiyun Gong, Ph.D.

Committee Member

David J. Cutler, Ph.D.

Committee Member

2022

Alzheimer's Disease Pathology Imputation and Risk Prediction Using Clinical Indices

By

Xizhu (Emilia) Liu


Jingjing Yang, Ph.D.

Adviser


An abstract of
a thesis submitted to the Faculty of Emory College of Arts and Sciences
of Emory University in partial fulfillment
of the requirements of the degree of
Bachelor of Science with Honors


Department of Quantitative Theory and Methods


2022

Abstract

Alzheimer's Disease Pathology Imputation and Risk Prediction Using Clinical Indices
By Xizhu (Emilia) Liu

Alzheimer's disease (AD) is a chronic progressive disorder that develops over years before manifesting impaired cognition, and early detection and intervention before the onset of noticeable AD symptoms might slow down the progression of cognitive decline. Since existing AD biomarkers are problematics and not widely available, this study aimed to develop models for imputing AD brain pathology and predict the risk of AD using common clinical indices. Data used in this study included clinical indices and postmortem pathology data contributed by 2000+ participants from two cohort studies who agreed on annual clinical visit and brain donation after death. In stage 1 of our study, we validated imputation models and chose the best-performing machine learning method: generalized linear regression model with elastic net regulation. In stage 2, we applied the imputation models to estimate baseline AD pathology using 57 clinical variables as predictors. In stage 3, we fitted Cox proportional hazard models and used the imputed pathology along with three demographic indices to predict the risk of cognitive impairment and AD dementia over years. Based on our data analysis results, imputed pathology was able to distinguish AD pathology-absent participants from AD pathology-present participants, and the clinical variables measured at baseline were effective predictors of baseline AD pathology. Moreover, imputed pathology along with three demographic indices were enough to make effective prediction on the risk of developing mild cognitive impairment or Alzheimer's disease dementia. If the leveraged clinical indices—common, affordable, and convenient to be measured—can be used as new biomarkers that substitute the existing but problematic ones, many more elderly people would be able to benefit from early detection, intervention and prognosis of their potential risk of developing AD dementia or cognitive impairment.

Alzheimer's Disease Brain Pathology Imputation and Risk Prediction Using Clinical Indices

By

Xizhu (Emilia) Liu

Jingjing Yang, Ph.D.

Adviser

A thesis submitted to the Faculty of Emory College of Arts and Sciences
of Emory University in partial fulfillment
of the requirements of the degree of
Bachelor of Science with Honors

Department of Quantitative Theory and Methods

2022

Acknowledgements

Table of Contents

**1. Introduction**

**1.1. Alzheimer's disease**

Alzheimer's disease (AD) is a chronic progressive disorder that develops over many years before

manifesting impaired cognition. In its earliest asymptomatic stage, Alzheimer's disease

pathology accumulates in adults with normal cognition and higher levels of AD are associated

with an increased risk of Alzheimer's dementia (Ballard et al.). While AD is a pathological

disease, dementia is a general term that summarizes a set of symptoms associated with

memory loss, limited social skills, impaired thinking abilities, or declined daily functioning, etc.

Accounting for approximately 70% of all dementia cases, AD is a type of and the most common

cause of dementia (Boyle et al.). As complex brain changes develop and cell damage

accumulates, dementia symptoms tend to occur and worsen over time (Yu et al.). Mild

cognitive impairment (MCI) causes cognitive declines that are not yet severe enough to affect

daily life or independent function. In terms of severity, MCI can be regarded as the precursor of

AD dementia—which is a severer case—but does not necessarily lead to AD dementia if the

progression is slow (Larson et al.).

Of the total U.S. population, around 6.5 million people aged 65 and above are currently living

with Alzheimer's dementia in the U.S. in 2022. Among those, nearly 75% are of age 75 or older.

About one in every nine people aged 65 and above suffers from Alzheimer's dementia.

According to recent studies, the population of age 65 and above with Alzheimer's dementia is

predicted to be as large as 12.7 million by year 2050. Because there could be nearly 20 years of

lagging from the start of pathological changes to the onset of AD symptoms, AD has been difficult to be predicted, diagnosed or treated in time (Alzheimer's Association).

Although AD is currently considered an irreversible disease, scientists and physicians suggest that early intervention could slow down the development of cognitive decline and the progression to the final stage of AD (Karikari et al.). Therefore, it is important to make predictions and recognize early signs before the onset of noticeable AD symptoms. Consequently, there has been an intense focus on identifying AD biomarkers that can estimate the burden of AD to facilitate identification of at-risk adults who might benefit from early interventions to reduce the accumulation of AD and prevent the development of cognitive impairments due to AD (Lan et al.).

However, existing AD biomarkers such as cerebrospinal fluid (Counts et al.) are not widely available due to their costs of advanced brain imaging and neuropsychological testing, invasiveness and difficultly to deploy at scale (Darst et al.). Therefore, it has been of our great interest to investigate and discover predictive and diagnostic methods that are more affordable, convenient, noninvasive, and easy to deploy at scale.

## 1.2. Motivations

Studies have so far based the biological understanding of AD on pathological changes in the brain. Brain pathology refers to the features, conditions, or typical behaviors of a disease that are reflected in the biological appearance or changes in the brain, including but not limited to the expression and accumulation of certain proteins. There are two neural proteins that signify

the presence of AD: an outside-neurons protein beta-amyloid and an inside-neurons protein tau (tangles) with twisted strands. These pathological changes are likely to be followed by the apoptosis of neurons and damage on brain tissues, and higher levels of accumulated pathology are associated with higher risk of AD. Besides beta-amyloid and tangles, there are two other pathology indices involved in this study: Global AD Pathology which is a quantitative summary of the expression of the two proteins beta-amyloid and tangles, and NIA-Reagan which is a dichotomous diagnosis criterion that defines and recodes the presence of AD pathology based on consensus recommendations for postmortem diagnosis of AD (Newell et al.). A high or intermediate likelihood of AD pathology is recoded as NIA-Reagan = 1, and a low or none likelihood of AD pathology is recoded as NIA-Reagan = 0 (Rush Alzheimer's Disease Center).

Indices of AD pathology necessary for a pathologic diagnosis of AD can only be obtained after death via autopsy, which puts significant challenge for estimating the burden of AD pathology in living adults during the course of AD. Prior studies have used a combination of brain imaging and clinical metrics to impute the burden of AD pathology, but brain imaging is expensive and not widely available (Counts et al.). A wide variety of analytic strategies and machine learning techniques (Hastie et al.) have been employed in diverse areas of aging research to impute or estimate missing data and data that are difficult to collect. These analytical approaches have been extended to estimate transcriptomic data that are generally hard to profile due to high cost and limited accessibility to some human tissues such as brain and kidney tissues, by using genetic genotype data that can be easily profiled from whole blood with a low cost as

predictors (Gamazon et al.). As a result, integrative analysis with both transcriptomic and genetic data is practically feasible.

Building on these prior studies to impute missing data due to limited accessibility, clinical-autopsy studies which collect both clinical indices prior to death and postmortem AD indices might be used to develop imputation models with postmortem samples. These imputation models can then be applied in living adults using clinical measures as predictors to estimate the burden AD pathology. If successful, this approach may yield a more widely available AD biomarker, since once an imputation model for AD has been validated, the model can then be applied to estimate AD in any sample of living adults with analogous clinical indices.

**1.3. Study design**

The main purpose of this study was to develop models with sufficient predictive ability to estimate the levels of AD pathology in living adults without postmortem autopsy, and to predict the risk for low-risk adults to develop MCI or even ADD over years. This was an observational study that collected data from both annual clinical visits for living adults who continued annual clinical visit, as well as postmortem brain autopsy for people after death. A total of 57 clinical indices measured in living adults were initially considered as predictors or covariates to be used in our models, which were expected to be possible substitutes for postmortem AD pathology data. Meanwhile, experts were able to obtain through autopsy the quantitative levels of the four AD pathology of our interest that were observed in those postmortem brains: beta-amyloid, tangles, Global AD Pathology, and NIA-Reagan.

Participants recruited for this research were community dwelling older persons enrolled in one of two ongoing cohort studies of aging and dementia, the Religious Orders Study [ROS, N = 1104 (48.3%)] and Rush Memory and Aging Project [MAP, N = 1183 (51.7%)]. Both studies were approved by the Institutional Review Board of Rush University Medical Center. Participants were enrolled without known dementia and agreed to annual clinical follow-up evaluations and autopsy at the time of death. Both studies share a large common core of testing batteries and uniform structured clinical evaluations by the same staff facilitating combined analysis. A written informed consent and an anatomical gift act were obtained from each participant.

## 1.4. Research goals and objectives

In this study, we aimed to (1) derive imputation models for AD brain pathology using clinical variables as predictors, (2) estimate the level of AD pathology in living adults, and (3) develop risk prediction models that identify adults at risk for AD and cognitive impairment using clinical variables alone. The current study used clinical and postmortem indices from more than 2000 decedents, who had participated in two community-based cohort studies and underwent brain autopsy at death, to train and validate imputation models that estimate the burden of AD pathology based on clinical measures alone. The estimated levels of AD pathology will be further used for creating models for AD risk prediction over years.

## 2. Methods

### 2.1. Clinical variables

The two main datasets that were used in this study included a longitudinal dataset and a cross-sectional dataset. The longitudinal dataset tracked each participant's annual clinical measurements from their baseline (or first) clinical visit to the last visit they were able to make. The cross-sectional dataset contained five groups of clinical indices across different measurement categories, composing a total of 57 individual clinical variables.

The five groups of clinical variables were considered as predictors for the burden of AD: i) Parsimonious AD risk factors including age, sex, education, cognitive assessment based on Mini Mental State Examination (MMSE) (Tombaugh et al.) and APOE E4 genotype; ii) Clinic variables and chronic health variables conditions such as blood pressure, depression, and cardiovascular disease, besides; iii) Medication usage variables such as antibiotics, lipid lowering medicine, and antidepressants; iv) Variables uniquely profiled by ROS/MAP such as global cognition test scores based on 17 cognitive tests, self-reported physical and social activities (Bennett et al.); v) Motor and sleep variables measuring patients such as gait speed, hand strength, global parkinsonian score (Buchman et al.), and 4 survey questions about sleep status (Park et al.). These motor variables have the potential to confound the associations of motor abilities with total daily physical activity by degrading motor capacity or affect an individual's propensity to engage in physical (Supplemental Table 1).

| Events | NCI | MCI | ADD | Row Total |
|---|---|---|---|---|
| Baseline NCI | 978 | 370 | 398 | 1746 |
| Baseline MCI | 196 | 245 | 100 | 541 |
| Column Total | 1174 | 615 | 498 | 2287 |

**Table 2. Sample sizes per event type for the ROS/MAP cohorts used for fitting prediction**

**models of the burden of AD pathology.**
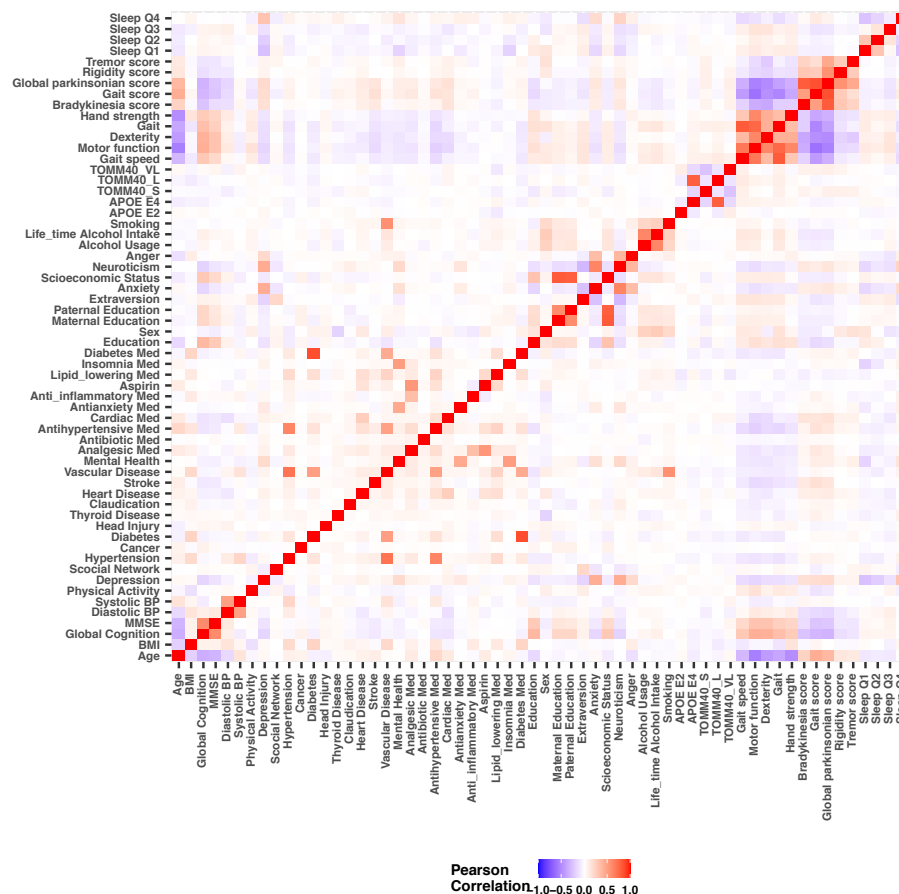


**Figure 1. Correlation heatmap of clinical variables as predictors of the burden of AD**

**Pathology.**

The above heatmap displays the magnitudes of correlations between every two clinical

variables, with each small grid representing the correlation of one pair of clinical variables. The

lighter the grid color, the weaker the correlation; the darker the grid color, the stronger the

correlation. Most grids displayed light or even white colors indicating that most of our clinical variables exhibited weak to extremely weak correlations with other covariates, so multicollinearity was not likely to be a problem in the regression models. Only a few exceptions were present in the heatmap: some composite motor variables displayed dark red or dark purple near the top-right of this heatmap and indicates high collinearity. Three composite motor variables thus needed to be excluded from the regression models later in this study.

## 2.2. Assessment of AD brain pathology

Brain autopsy followed a standard protocol (Schneider et al.). Neuropathologic evaluations, blinded to clinical data, assessed the burden of four measures of AD. A modified Bielschowsky silver stain was used to visualize neuritic plaques, diffuse plaques, and neurofibrillary tangles in five cortical areas (hippocampus, entorhinal, midfrontal, middle temporal, and inferior parietal). Neuritic and diffuse plaques, and neurofibrillary tangles were counted in the region that appeared to have the maximum density of each pathology as previously described. A standardized score was created for each neuropathology in each region by dividing the raw count by the standard deviation of the mean for the same neuropathology in the same region. This standardization procedure puts the pathologic indices on a relatively common scale. They were averaged to create a composite measure as previously described (Bennett et al.). The National Institute on Aging-Reagan criteria were used with intermediate and high likelihood cases indicating a pathologic diagnosis of AD (The National Institute on Aging).

In addition, a summary global AD pathology score was made based on the greatest density of neuritic plaques, diffuse plaques, and neurofibrillary tangles in one mm2 (Bennett et al.). Amyloid-β was labeled with an N-terminus–directed monoclonal antibody (10D5; Elan, Dublin, Ireland; 1:1,000). Immunohistochemistry was performed as previously described using diaminobenzidine as the reporter, with 2.5% nickel sulfate to enhance immunoreaction product contrast. PHFtau was labeled with an antibody specific for phosphorylated tau (AT8; Innogenetics, San Ramon, CA; 1:1,000). Amyloid-β (Aβ) load and tangles (tau tangles) were quantified in 8 brain regions (anterior cingulate cortex, superior frontal cortex, mid frontal cortex, inferior temporal cortex, hippocampus, entorhinal cortex, angular gyrus/supramarginal gyrus, and calcarine cortex). Overall Aβ load was calculated through averaging mean percent area of Aβ deposition per region, across multiple brain regions. Tangles densities were derived by averaging tangles densities across corresponding brain regions.

## 2.3. Data preparation

Since the analysis of longitudinal clinical data relies on completeness of data in each clinical variable in each year, missing values due to participants' occasional absence of clinical visits or some of the clinical examinations must be imputed and filled. For each participant, missing values after baseline visit were filled by the values of the nearest previous visit, while missing values at baseline visit were filled by the values of the nearest next visit. Besides imputing missing values, we prepared the data by standardizing all variables, including both the clinical variables and the four AD pathology. The variables were standardized at mean of 0 and standard deviation of 1.

## 2.4. Analytic approach

### 2.4.1. Stage 1: Train and validate imputation models for AD pathology

The first stage of this research was model validation. We applied cross validation and compared among several potential approaches and imputation models. We used near-death (last-visit) data from MAP participants to fit the training models and near-death data from ROS participants to test the fitted models. In the training models, the last-visit clinical variables in MAP were used as predictors and postmortem pathology data in MAP were used as responses. In the testing models, the last-visit clinical variables in ROS were used as predictors to estimate the postmortem pathology data in ROS (Figure 1).



**Figure 2. Stage 1 flowchart: model validation by training on near-death MAP cohort and testing on near-death ROS cohort**

To choose the optimal imputation model that returns the highest prediction accuracy, we experimented and compared four machine learning methods: Gradient Boost Machine (GBM), Support Vector Machine (SVM), Random Forest, and Generalized Linear regression Model with Elastic-Net penalty (GLM-EN). Among these four, we compared the prediction $R^2$ and area under curve (AUC) values corresponding to each of the four brain pathologies yielded by each machine learning method. The comparison would give us the optimal method for fitting imputation models. We used last visit data in MAP participants to fit the training models and

last visit data in ROS participants to test the fitted models. Each model was fitted for one of the four brain pathologies: Tangles, Global AD Pathology, Beta-amyloid, and NIA-Reagan.

In each GBM model, we isolated the target brain pathology and fitted that pathology on all other clinical variables except the other three brain pathologies. For the continuous pathologies Tangles, Global AD Pathology, and Beta-amyloid, we specified Gaussian distributions in the model fitting process. For the binary pathology NIA-Reagan, we specified Bernoulli distributions and found a threshold that produced ~80% sensitivity on predicted probabilities. We performed 10-fold cross validation at the optimum number of trees/iterations in each model and then obtained the prediction $R^2$ and AUC values.

Within each SVM model, we compared among four different types of kernel used for training and predicting, including linear, polynomial, radial, and sigmoid. Linear kernel produced the highest AUC values across all models and was thus chosen to be the kernel type we used in SVM. The cost of constraints violation was set to its default value of 1, which gave us the best prediction $R^2$ and AUC values.

For GLM-EN, we used the "cv.glmnet" function to select the alpha and lambda levels which were later used in our model-fitting process with the "glmnet" and "predict" functions. We produced prediction $R^2$ and AUC values for each model, as well as box plots with p-values obtained from two-sample hypothesis T tests. We chose the best-performing model and continued applying it to the pathology imputation process in stage 2, while the other three models were excluded from the rest of this study.

### 2.4.2. Stage 2: Estimate AD pathology at study baseline

GLM-EN was the best-performing method and was chosen as the approach for pathology imputation. Since AD brain pathologies are not measurable in living adults but can only be measured in postmortem brains, we trained imputation models with near-death brain pathology data and imputed for brain pathologies at baseline visit. In the training models, last-

visit clinical variables in all (both MAP and ROS) participants were used as predictors, and postmortem pathology data in all participants were used as responses. Applying these fitted models, we then used baseline clinical variables in all participants as predictors to estimate all participants' pathology levels at baseline, which cannot be obtained otherwise through autopsy in living adults (Figure 3).
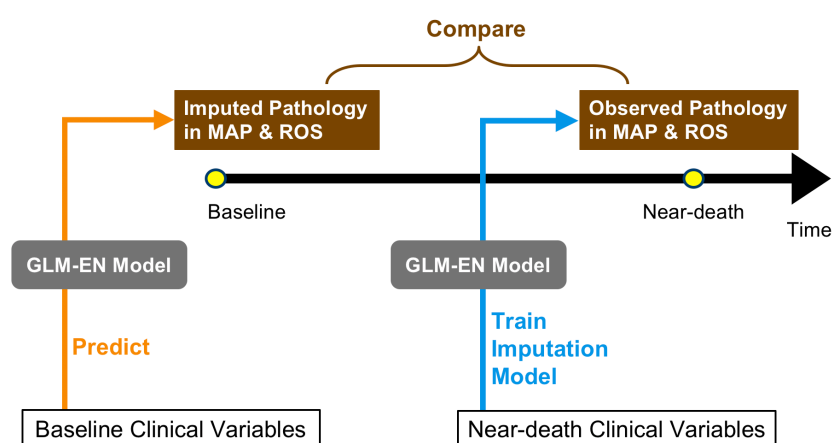


**Figure 3. Stage 2 flowchart: AD pathology imputation by using GLM-EN models to train on near-death MAP & ROS cohort and to predict pathology level in MAP & ROS at baseline.**

For each model, we used 57 clinical metrics as predictors for the imputed pathology. The performance of each imputation model was measured by Pearson correlations, prediction R-squared, 95% confidence intervals, and p-values. After each of the four brain pathologies was imputed for their baseline values, the newly estimated pathology data at baseline were stored as new vectors or variables, prepared for AD risk prediction later in the 3rd stage of the study.

### 2.4.3. Stage 3: Predict AD risk by Cox proportional hazard model

With AD pathology data imputed at baseline, we were further interested in whether the estimated level of baseline AD pathology was associated with incident MCI, incident AD dementia (ADD) and the risk of pathologic AD based on NIA-Reagan at autopsy. The main goal of stage 3 of this study was to examine whether the estimated levels of brain pathology (from

stage 1 and 2) together with some selected clinical variables could effectively and accurately predict the occurrence of incident ADD or MCI, and whether the prediction based on models involving both imputed pathology and clinical metrics would be improved as compared to models using only clinical indices (Figure 3).



**Figure 4. Stage 3 flowchart: risk prediction models by training on baseline clinical variables and imputed pathology and predicting occurrence of MCI and/or ADD events in year 3 and 5**

During follow-up years after the participants' first visits at baseline (which was considered as year 0), some participants were identified with occurrence of Alzheimer's dementia (ADD) or mild cognitive impairment (MCI) through annual cognitive status diagnosis. The year of first diagnosis of ADD or MCI was considered as the time when the event/incident occurs. For living participants, the last time of their annual visit were considered as the right censored time, while for dead participants, their last visit records were considered as their year of death.

We fitted respective Cox proportional hazard risk prediction models for the events of ADD and MCI while accounting for the competing risk of death. We evaluated the accuracy of the risk prediction models by constructing receiver operating characteristic (ROC) curves and calculating the area under curve (AUC) with respect to the participants' observed ADD or MCI status in the given year. The covariates that were included in each of our final Cox regression models were chosen through backward selection, which was implemented during the training of the corresponding Cox models that involved more than one covariate.

## 3. Results

### 3.1. Stage 1: Train and validate imputation models for AD pathology

We incorporated 16 individual models in the model validation process, applying each of the four potential machine learning methods on each of the four AD pathology. As shown in Table 3, the predictive ability of each model was represented by either a value of prediction $R^2$ or a value of area under ROC curve (AUC). For the three continuous variables Beta-Amyloid, Tangles, and Global AD Pathology, each prediction $R^2$ value was obtained by squaring the correlation between the observed levels and the estimated levels of near-death AD pathology in the ROS testing cohort. For the binary variable NIA-Reagan, each model's predictive ability was measured by the area under ROC curve, reflecting the performance the model in distinguishing between two diagnostic groups.

| Machine learning method | Prediction R² | | | AUC |
| --- | --- | --- | --- | --- |
| | Beta-Amyloid | Tangles | Global AD Pathology | NIA Reagan |
| GLM-EN | 0.170 | 0.332 | 0.248 | 0.761 |
| GBM | 0.144 | 0.305 | 0.221 | 0.740 |
| SVM | 0.113 | 0.270 | 0.190 | 0.632 |
| Random Forest | 0.152 | 0.305 | 0.218 | 0.731 |

**Table 3. Comparison of machine learning method performance: results of prediction R² and AUC obtained from cross validation models in stage 1.**

Comparing among the four machine learning methods, we found that GLM-EN yielded the highest prediction R² and AUC in all of the four AD pathology. Therefore, GLM-EN was chosen as the machine learning method for imputing the AD pathology levels at baseline in stage 2.
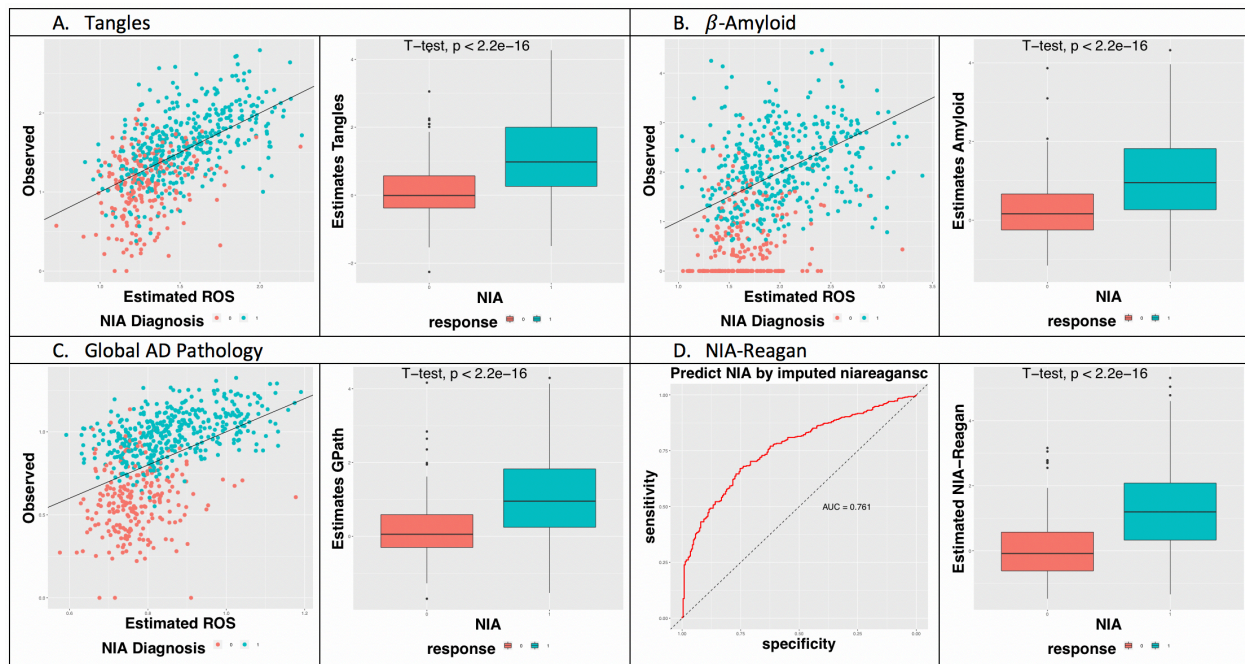


**Figure 5. Pathology estimation results using GLM-EN in ROS testing samples at death and their discrimination with respect to profiled NIA-Reagan at death for Tangles (A), Beta-Amyloid (B), Global AD pathology (C), and NIA-Reagan (D).**

Figure 5 is a visualization of the performance of the GLM-EN imputation models corresponding to each of the four AD pathology. For the three continuous variables (A, B, C), observed pathology levels were plotted against estimated pathology levels. Each dot on the scatterplots represented an individual participant, belonging to either the pathology-absent group (red, NIA = 0) or the pathology-present group (blue, NIA = 1). For the binary variable (D), a scatterplot was replaced by a ROC curve with calculated AUC. In addition, a boxplot was plotted for each of the four AD pathology, showing how the mean level of a pathology differs between the pathology-absent group (red, NIA = 0) and the pathology-present group (blue, NIA = 1). We performed two-sample t tests and obtained the corresponding p-value for each boxplot, finding that the p-values are statistically significant in all of the four pathology. Therefore, our imputation models were able to significantly distinguish between the two diagnostic groups based on the estimated AD pathology levels.

**3.2. Stage 2: Estimate AD pathology at study baseline**

In stage 2 of this study, we imputed the AD pathology by fitting models on all near-death data and estimating the burden of AD pathology at baseline. Figure 6 summarized the effect sizes of some of the covariates that were selected by the GLM-EN model based on their significance in each pathology imputation model, with beta values on x-axis showing the covariate effect sizes of the standardized predictors. Longer bars in the plots represented greater magnitude and

stronger partial effect of a predictor on the estimated pathology levels, regardless in the
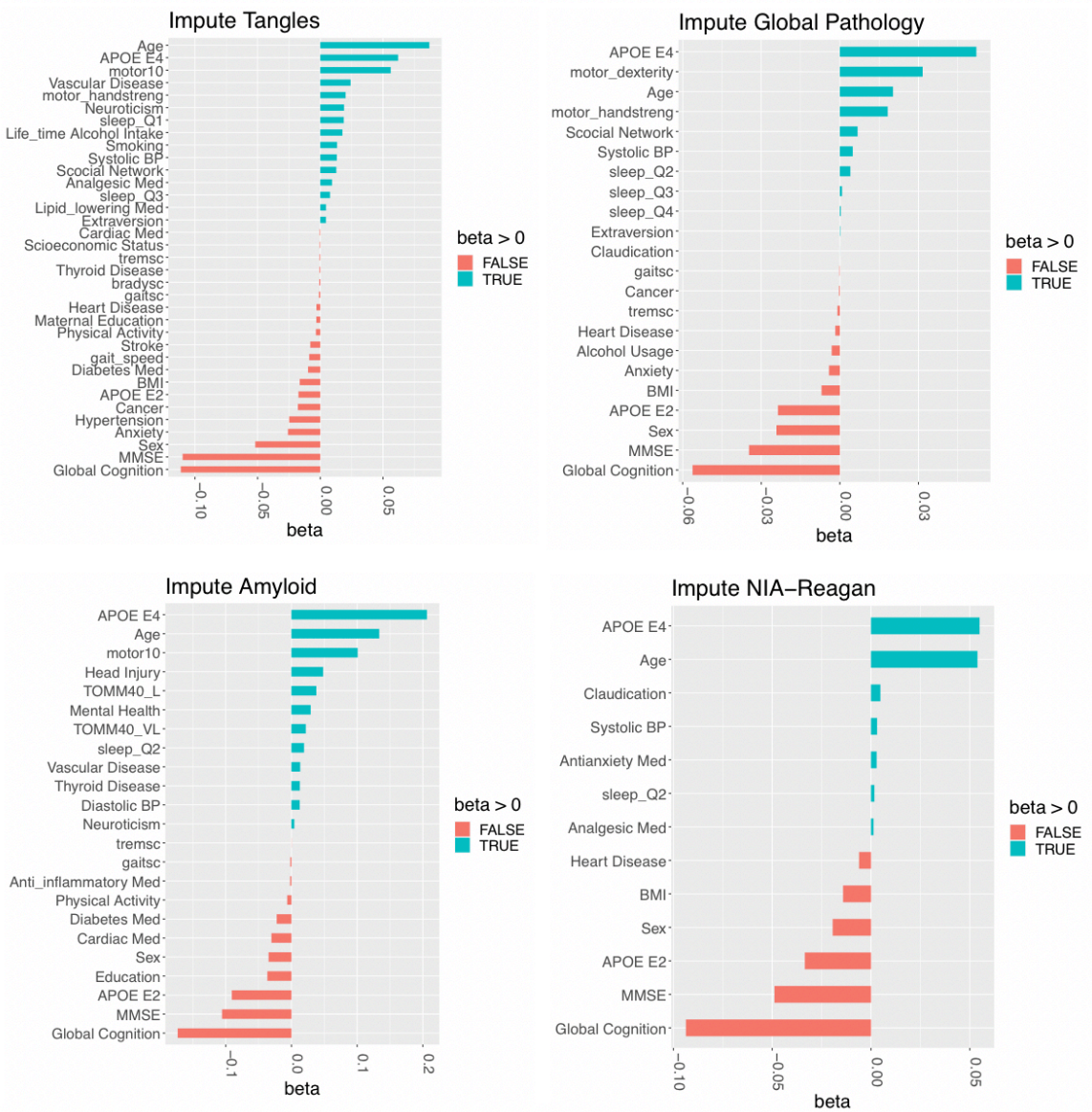
positive or negative direction.



**Figure 6. Important covariates selected by GLM-EN in the prediction models for the burden of AD pathology.**

A few predictors were commonly selected across the four pathology. For instance, age and the

APOE E4 gene were considered as important predictors that positively contribute to the

increase in AD pathology levels. Older ages might be associated with higher burden of AD pathology, and the expression level of APOE E4—a gene that increases the risk of developing AD—might also lead to higher estimated pathology levels. Moreover, global cognition score, MMSE (Mini Mental State Exam), and the APOE E2 gene were selected as influential covariates that negatively contribute to the burden of AD pathology. Higher scores in global cognition test and MMSE might indicate better thinking and memory ability and thus tended to decrease the estimated level of AD pathology. The expression level of APOE E2—a gene that is associated with low risk of developing AD—might also lead to lower estimated pathology levels.

|  |  | Beta-Amyloid | Tangles | Global AD Pathology | NIA-Reagan |
|---|---|---|---|---|---|
| **Training Results** | **Pearson Correlation or AUC** | 0.470 | 0.627 | 0.542 | 0.832 |
| **Prediction Results** | **Pearson Correlation or AUC** | 0.319 | 0.394 | 0.412 | 0.744 |
|  | **Correlation test or t test P-value** | 4.074e-27 | 5.539e-42 | 1.415e-47 | 7.956e-34 |

**Table 4. Training and testing results obtained from GLM-EN imputation models in stage 2.**

After using GLM-EN models to impute the baseline values of the four AD pathology, we measured the predictive ability of each model by computing either Pearson correlation or AUC, along with the corresponding p-values obtained from correlation tests or t tests. Each correlation compared imputed pathology levels at baseline with observed pathology levels profiled at death. Comparison across different time stages in the progression of AD was based on the fact that baseline AD pathology follows a certain pattern to develop into near-death AD pathology, and a significant correlation might suggest that the imputation model was able to

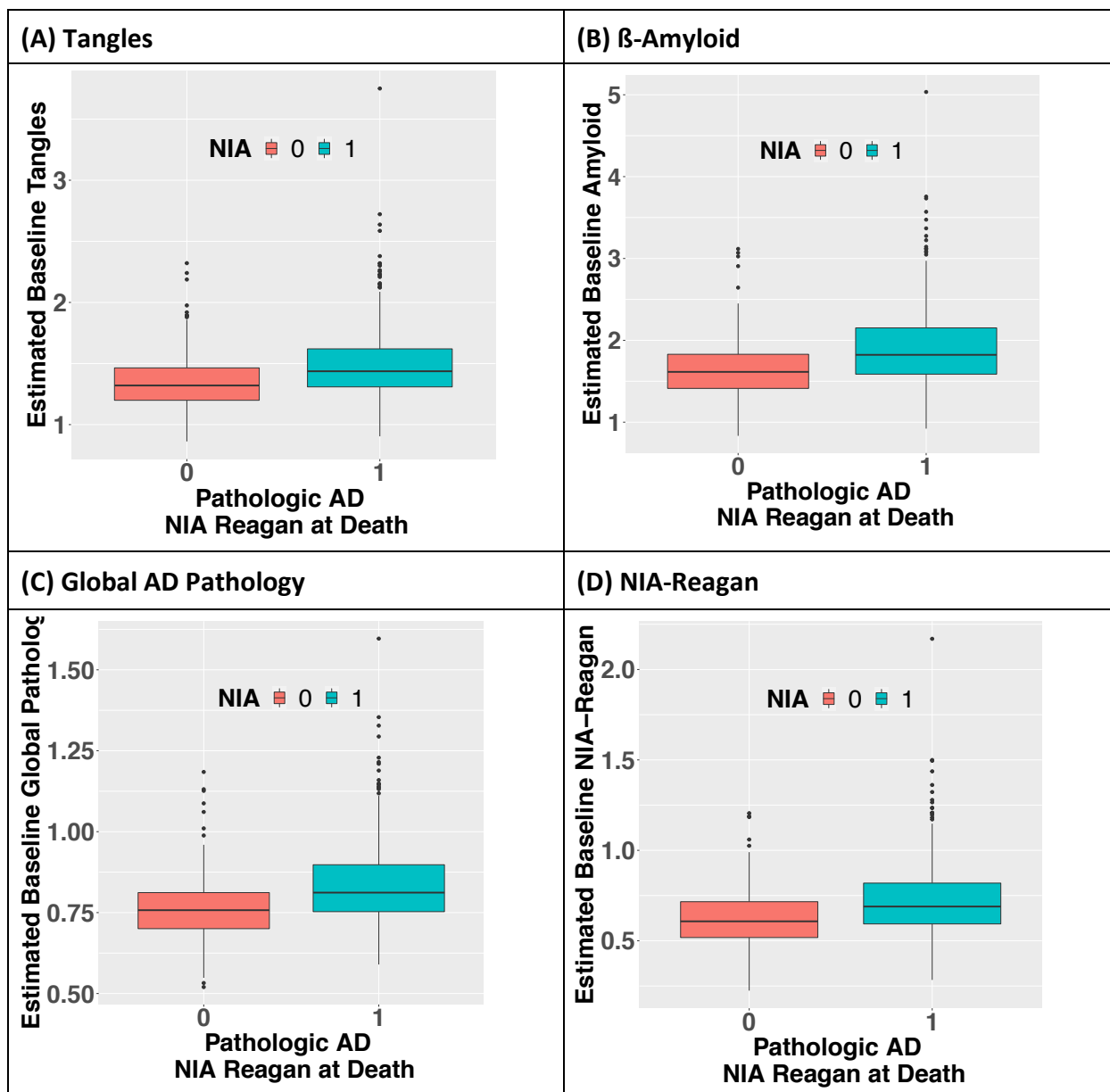capture that pattern and thus make meaningful estimation on the burden of AD pathology at baseline.



**Figure 7. Imputation performance at baseline with respect to NIA-Reagan at death (stage 2).**

**Boxplots of predicted AD pathology at baseline with respect to profiled NIA-Reagan at death.**

To determine whether the mean in each estimated baseline AD pathology significantly differs between the two diagnostic groups that were categorized based on NIA-Reagan profiled at death, we plotted a boxplot for each AD pathology and computed the t test p-values. All of the four p-values were statistically significant, meaning that the imputed pathology was able to distinguish pathology-absent participants from pathology-present participants, and that the clinical variables measured at baseline were effective predictors of baseline AD pathology.

### 3.3. Stage 3: Predict AD risk over years

After imputing the burden of AD pathology at baseline in stage 2, the stage 3 of this study developed Cox proportional hazard regression models to predict the risk of MCI and ADD events occurrence at year 3 and year 5 after participants' baseline visits. Each Cox model involved four predictors: age at baseline visit, sex, education, and one of the four imputed baseline pathology. Shown in Figure 8, we separated the occurrence of MCI and ADD events into three different scenarios: (A) participants with either NCI (no cognitive impairment) or MCI at baseline developing ADD in year 3 or 5, (B) participants with NCI at baseline developing ADD in year 3 or 5, and (C) participants with NCI at baseline developing MCI in year 3 or 5. Each cell with two ROC curves—the red curve representing year 3 and the blue curve representing year 5—is a visualization of the ability of the Cox models to correctly distinguish and diagnose participants with MCI or ADD.

**Figure 8. Risk prediction results of MCI and ADD events in year 3 and 5 by Cox proportional hazard regression models, using age, sex, education and imputed AD pathology as predictors.**

Each cell is also a comparison between the year-3 curve and the year-5 curve. The curve with

larger area under the ROC curve was considered a better prediction, and the prediction model

corresponding to that curve, either a year-3 model or a year-5 model, is considered a more

accurate model. As shown in Figure 8, in every cell across the three scenarios (rows) and the four models (columns), the red curve had a larger AUC compared to the blue curve. This might indicate that the Cox models had better predictive ability when predicting for the risk of MCI and ADD occurrence three years after baseline than five years after baseline. The risk prediction result would be more reliable for years that are closer to the year of baseline visit.

## 4. Conclusions

Leveraging novel clinical and postmortem AD measures from the two cohort studies, we trained and validated imputation models for AD pathology that estimated AD pathology at study baseline years before death. The imputation models worked by mathematically explaining the variation of observed AD pathology in autopsied adults by their equivalence, in clinical indices. Individuals with higher levels of estimated pathology at baseline were at higher risk of developing AD and of a postmortem diagnosis of AD. The imputation models learned the predictive information of postmortem AD from clinical indices of decedents undergoing autopsy, and were able to estimate the burden of AD pathology levels using clinical indices alone. Once an imputation model for AD pathology has been validated, this approach might yield a more widely available AD biomarker, and the imputation model could be applied to estimate AD in any sample of living adults with analogous clinical indices. Since the information provided by the clinical indices were already incorporated and reflected in the estimated levels of AD pathology, any imputed pathology with the addition of merely three demographic indices were enough to make effective prediction on the risk of developing mild cognitive impairment or Alzheimer's disease dementia.

Moreover, if the leveraged clinical indices—common, affordable, and convenient to be measured—can be used as new biomarkers that substitute the existing but problematic ones, many more elderly people would be able to benefit from early detection, intervention and prognosis of their potential risk of developing AD dementia or cognitive impairment. Given that no cure has been invented to completely overcome Alzheimer's, predictive measures would play a significant role in lowering the risk of AD or slowing down the progression of the disease or related symptoms, which might bring elderly people higher life quality and even longer life span.

**Appendix**

| Group | Variable | Mean (SD) or N (%) |
|---|---|---|
| Clinical and chronic health variables | BMI (9.09, 62.91) | 27.52 (5.50) |
| | Diastolic blood pressure (40, 122.5) | 74.47 (11.90) |
| | Hypertension blood pressure (83, 215.5) | 134.13 (18.29) |
| | Depression score (0, 9) | 0.94 (1.48) |
| | Hypertension (N) | 1142 (49.93) |
| | Cancer (N) | 719 (31.43) |
| | Diabetes (N) | 277 (12.11) |
| | Head injury (N) | 148 (6.47) |
| | Thyroid disease (N) | 440 (19.23) |
| | Claudication (N) | 138 (6.03) |
| | Heart disease (N) | 208 (9.09) |
| | Stroke (N) | 156 (6.82) |
| | Cardiovascular disease history counts (0, 3) | 0.93 (0.78) |
| | Alcohol usage in the past year, grams per day (0, 116.55) | 4.56 (11.14) |
| | Alcohol usage when drank most in lifetime (0, 6) | 0.46 (0.98) |
| | Smoking (never smoked 0, former smoker 1, current smoker 2) | 0.33 (0.51) |
| Medication usage | Mental health (N) | 529 (23.13) |
| | Analgesic (N) | 1677 (73.32) |
| | Antibiotic (N) | 158 (6.90) |
| | Anti-hypertensive (N) | 1394 (60.95) |
| | Cardiac (N) | 228 (9.96) |
| | Anti-anxiety (N) | 133 (5.81) |
| | Anti-inflammatory (N) | 561 (24.52) |
| | Aspirin (N) | 992 (43.37) |
| | Lipid lowering (N) | 740 (32.35) |
| | Insomnia (N) | 171 (7.47) |
| | Diabetes (N) | 203 (8.87) |

**Supplemental Table 1.** Baseline characteristics of clinical variables by category

## References

Alzheimer's Association. "2022 Alzheimer's Disease Facts and Figures." Alzheimers Dement 2022, 2022.

Ballard C, Gauthier S, Corbett A, Brayne C, Aarsland D, Jones E. Alzheimer's disease. Lancet 2011;377:1019-1031.

Bennett DA, Buchman AS, Boyle PA, Barnes LL, Wilson RS, Schneider JA. Religious Orders Study and Rush Memory and Aging Project. J Alzheimers Dis 2018;64:S161-S189.

Bennett DA, Wilson RS, Boyle PA, Buchman AS, Schneider JA. Relation of neuropathology to cognition in persons without cognitive impairment. Ann Neurol 2012;72:599-609.

Bennett DA, Wilson RS, Schneider JA, et al. Education modifies the relation of AD pathology to level of cognitive function in older persons. Neurology 2003;60:1909-1915.

Boyle PA, Yu L, Leurgans SE, et al. Attributable risk of Alzheimer's dementia attributed to age-related neuropathologies. Ann Neurol 2019;85:114-124.

Buchman AS, Yu L, Oveisgharan S, et al. Cortical proteins may provide motor resilience in older adults. Scientific Reports 2021;11:11311.

Consensus recommendations for the postmortem diagnosis of Alzheimer's disease. The National Institute on Aging, and Reagan Institute Working Group on Diagnostic Criteria for the Neuropathological Assessment of Alzheimer's Disease. Neurobiol Aging 1997;18:S1-2.

Counts SE, Ikonomovic MD, Mercado N, Vega IE, Mufson EJ. Biomarkers for the Early Detection and Progression of Alzheimer's Disease. Neurotherapeutics 2017;14:35-53.

Darst BF, Lu Q, Johnson SC, Engelman CD. Integrated analysis of genomics, longitudinal metabolomics, and Alzheimer's risk factors among 1,111 cohort participants. Genet Epidemiol 2019;43:657-674.

Gamazon ER, Wheeler HE, Shah KP, et al. A gene-based association method for mapping traits using reference transcriptome data. Nature genetics 2015;47:1091-1098.

Hastie T, Tibshirani R, Friedman JS. The Elements of Statistical Learning, 2 ed: Springer-Verlag New York, 2009.

Karikari TK, Benedet AL, Ashton NJ, et al. Diagnostic performance and prediction of clinical progression of plasma phospho-tau181 in the Alzheimer's Disease Neuroimaging Initiative. Mol Psychiatry 2021;26:429-442.

Lan MJ, Ogden RT, Kumar D, et al. Utility of Molecular and Structural Brain Imaging to Predict Progression from Mild Cognitive Impairment to Dementia. J Alzheimers Dis 2017;60:939-947.

Larson, E B, et al. "Cognitive Impairment: Dementia and Alzheimer's Disease." Annual Review of Public Health, vol. 13, no. 1, May 1992, pp. 431–449, pubmed.ncbi.nlm.nih.gov/1599598/, 10.1146/annurev.pu.13.050192.002243. Accessed 10 Apr. 2021.

Newell, Kathy L., et al. "Application of the National Institute on Aging (NIA)-Reagan Institute Criteria for the Neuropathological Diagnosis of Alzheimer Disease." Journal of Neuropathology and Experimental Neurology, vol. 58, no. 11, Nov. 1999, pp. 1147–1155, pubmed.ncbi.nlm.nih.gov/10560657/, 10.1097/00005072-199911000-00004. Accessed 10 Apr. 2022.

Park M, Buchman AS, Lim AS, Leurgans SE, Bennett DA. Sleep complaints and incident disability in a community-based cohort study of older persons. Am J Geriatr Psychiatry 2014;22:718-726.

Rush Alzheimer's Disease Center. "Variable Details | RADC." Www.radc.rush.edu, RADC Research Resource Sharing Hub, www.radc.rush.edu/docs/var/detail.htm?category=Pathology&subcategory=Alzheimer%27s+disease&variable=niareagansc. Accessed 10 Apr. 2022.

Schneider JA, Arvanitakis Z, Leurgans SE, Bennett DA. The neuropathology of probable Alzheimer disease and mild cognitive impairment. Annals of neurology 2009;66:200-208.

Tombaugh TN, McIntyre NJ. The mini-mental state examination: a comprehensive review. J Am Geriatr Soc 1992;40:922-935.

Yu L, Wang T, Wilson RS, et al. Common age-related neuropathologies and yearly variability in cognition. Annals of clinical and translational neurology 2019;6:2140-2149.