

IDENTIFYING RISK FACTORS FOR COGNITIVE DECLINE
USING STATISTICAL LEARNING TECHNIQUES AND
FUNCTIONAL DATA ANALYSIS

A thesis submitted to the
College of Graduate and Postdoctoral Studies
in partial fulfillment of the requirements
for the degree of Master of Science
in the Department of Mathematics and Statistics
University of Saskatchewan
Saskatoon

By
Hao Hu

©Hao Hu, 2022. All rights reserved.

Unless otherwise noted, copyright of the material in this thesis belongs to
the author.

Permission to Use

In presenting this thesis in partial fulfillment of the requirements for a Postgraduate degree from the University of Saskatchewan, I agree that the Libraries of this University may make it freely available for inspection. I further agree that permission for copying of this thesis in any manner, in whole or in part, for scholarly purposes may be granted by the professor or professors who supervised my thesis work or, in their absence, by the Head of the Department or the Dean of the College in which my thesis work was done. It is understood that any copying or publication or use of this thesis or parts thereof for financial gain shall not be allowed without my written permission. It is also understood that due recognition shall be given to me and to the University of Saskatchewan in any scholarly use which may be made of any material in my thesis.

Disclaimer

Reference in this thesis to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement, recommendation, or favoring by the University of Saskatchewan. The views and opinions of the author expressed herein do not state or reflect those of the University of Saskatchewan, and shall not be used for advertising or product endorsement purposes.

Requests for permission to copy or to make other uses of materials in this thesis in whole or part should be addressed to:

Head of the Department of Mathematics and Statistics
140 McLean Hall
University of Saskatchewan
Saskatoon, Saskatchewan S7N 5C9 Canada

OR

Dean
College of Graduate and Postdoctoral Studies
University of Saskatchewan
116 Thorvaldson Building, 110 Science Place
Saskatoon, Saskatchewan S7N 5C9 Canada

Abstract

Background: Numerous studies have shown that older adults' cognitive abilities are age-related and likely to decline at a certain age. Based on these indications, this work uses functional principal component analysis (FPCA) to explore changes in cognitive function with age. This study aims to describe the longitudinal cognitive function of elderly trajectory patterns between 65 and 80 years of baseline age using FPCA and identify risk factors for cognitive decline using machine learning algorithms.

Methods: We used FPCA to extract the overall pattern change and use elastic-net, decision tree, and random forest models to find risk factors. In a sample of elderly ($n = 944$) with 6608 measurements (7 waves) from the Survey of Health, Aging, and Retirement in Europe (SHARE), by using age at interview as time, longitudinal cognitive function trajectory patterns for the elderly were extracted using FPCA. Zou and Hastie (2005) proposed the elastic net regression method, which effectively implements feature selection by setting coefficients of non-significant variables to zero [56]. Random forest is a tree-based machine learning algorithm that harnesses the power of multiple decision trees to make decisions. Random forests combine the outputs of individual decision trees to generate the final result. [15]. We modelled the first two functional principal components (FPCs) with these machine learning algorithms and used the selected covariates in the baseline wave to identify risk factors.

Results and Conclusions: We have obtained four FPCs explained by 78.0, 14.2, 6.7 and 1.1 % of the variation respectively for the cognitive function. The mean function of FPCA shows that cognitive decline is generally divided into two stages (early decline and late decline). By analyzing and comparing a set of models at the national level, the cognitive function of each country is slightly different. Older people in Italy and Spain have significantly lower cognitive abilities. The Predictive R^2 of FPC1 and FPC2 is around 0.5 and 0.2 with covariate delayed recall score (DRS) and reduced to 0.4 and 0.1 without covariate DRS. From the individual point of view, many risk factors are modifiable and can be prevented in advance. Our results show that immediate recall score, education level, country, numeracy score, reading score, and household income are associated with cognitive patterns in the elderly.

Acknowledgements

First of all, I would like to express to my supervisor Dr. Li Xing and Longhai Li, that without their continuous guidance and encouragement, I could not complete my thesis. It is my great honour to work with them. I'm also grateful to Dr. Steven Rayan, the Graduate Chair in the Department of Mathematics & Statistics and Prof. Junxin Liu, my committee member, for academic support during my master's program. Besides, this thesis is supported by the University of Saskatchewan (USask) through a USask Internal Research Grant and the USask College of Arts and Science. This thesis was partially enabled by support provided by Compute Canada (www.computeCanada.ca). The data used in this paper is from SHARE Waves 1 [6], 2 [7], 4 [8], 5 [9], 6 [10], 7 [11], and 8 [12].

Last, I am thankful to my parents, who have loved me and supported me all the time.

Contents

Permission to Use	i
Abstract	ii
Acknowledgements	iii
Contents	iv
List of Tables	vi
List of Figures	vii
List of Abbreviations	viii
1 Introduction	1
2 The data	5
2.1 Source of data	5
2.2 Primary Outcome	5
2.3 Covariates	6
2.4 Study sample	9
3 Methods	10
3.1 Sparse Functional principal component analysis	10
3.2 Statistical learning	12
3.2.1 Shrinkage Models	12
3.2.2 Decision tree	14
3.2.3 Random forest	16
3.2.4 Predictive R^2 with CV	18
3.2.5 Cross Validation	18
4 Result	21
4.1 FPCA	21
4.2 Comparison of FPC scores by countries	22
4.3 Predictive R^2 for models	25
4.4 Identification of Risk factors for cognitive decline	25
5 Conclusion	31
References	32
Appendix A Plots for predicting the first two FPC scores from the categorical variables	35
Appendix B R Code	37
B.1 Data Management for Wave8	37
B.2 Load data and data clearing	37
B.3 Load data and data clearing	45
B.4 FPCA	47
B.5 PC score V.s. nine countries	48
B.6 data management before model fitting	50
B.7 Model fitting	50

B.7.1	Elastic Net	50
B.7.2	Decision Tree model with discrete variables	52
B.7.3	Decision Tree model with whole selected variables	52
B.7.4	Random Forest with whole selected variables	53
B.8	Predictive R^2 for different models	53
B.8.1	Elastic-net	53
B.8.2	Lasso	54
B.8.3	Random Forest	55
B.8.4	Decision Tree	56

List of Tables

2.1	Descriptive Statistics Across Waves (discrete variable)	6
2.2	Descriptive Statistics Across Waves (continues variable)	9
4.1	Euro Health Consumer Index 2015	25
4.2	Predictive R^2 for machine learning models	25

List of Figures

3.1	Decision Tree	15
3.2	Bootstrapping	17
3.3	5-fold Cross-Validation	19
3.4	Leave One Out Cross-Validation	19
4.1	Summary Of FPCA Result	21
4.2	Overall Cognitive Function Trajectory Across Different European Countries Based On FPC1. (Mean function + Median FPC1 of each country)	23
4.3	Overall Cognitive Function Trajectory Across Different European Countries Based On FPC1 and FPC2. (Mean function + Median FPC1 + Median FPC2 of each country)	23
4.4	Bloxplot Of FPC1 Across Different European Countries	24
4.5	Bloxplot Of FPC2 Across Different European Countries	24
4.6	Covariate Effects on the First FPC Scores by Elastic-Net, Note all the covariates are in the baseline wave 1.	27
4.7	Covariate Effects on the Second FPC Scores by Elastic-Net	27
4.8	Pruned Tree based method for predicting FPC1 from all the variables	29
4.9	Pruned Tree based method for predicting FPC2 from all the variables	29
4.10	Top 30 ranked variables identified by Random Forest based on the Two FPC Scores	30
A.1	Pruned Tree based method for predicting FPC1 from the categorical variables	35
A.2	Pruned Tree based method for predicting FPC2 from the categorical variables	36

List of Abbreviations

SHARE	Survey of Health, Aging, and Retirement in Europe
PCA	Principal component analysis
FPCA	Functional principal component analysis
BMI	Body mass index
IQR	Interquartile Range=Upper Qualitle-Lower Qualitle
USask	University of Saskatchewan
AD	Alzheimer’s disease
PD	Parkinson’s Disease
FPCs	Functional principal components
LHC	Life History Calendar
DRS	Delayed recall score
IRS	Immediate recall Score
FCS	Conditional specification method
E	Expectation
Var	Variance
MLTC	median longitudinal trajectory curves
CP	Complexity Parameter
FDA	Functional Data Analysis
EFs	Orthonormal eigenfunctions
i.i.d.	Independent and identically distributed
%IncMSE	Mean Decrease Accuracy
RSS	Residual Sum of Square
PRESS	Predicted Residual Sum of Squares
EHCI	Euro Health Consumer Index
PRIS	Patient rights and information score
WTTS	Waiting times for treatment score
OS	Outcomes score
RRSS	Range and reach of services score
PS	Pharmaceuticals score
RMSE	Root Mean Square Error
MAE	Mean Absolute Error

1 Introduction

Cognitive functions refer to various mental abilities, which are more related to how people learn, think, remember, pay attention, make decisions and solve problems [31]. Cognitive functions have been distinguished by how they change throughout adulthood, namely cognitive mechanics and cognitive pragmatics. Cognitive mechanics are constructed by information processing rate, working memory and inhibition (ability to automatically inhibit goal-irrelevant information). The information processing rate is most sensitive to age differences. In this work, we focus on cognitive mechanics, which is typically decline with age and accelerate in old age. In contrast to the mechanics, the pragmatics of cognition are verbal and numerical abilities that rely on the accumulation of lifespan development and knowledge-based forms of intelligence, which are well maintained during adulthood and will not decline obviously in advanced age [2].

Cognitive abilities are critical to living independently as people age, such as whether a person can live independently and drive safely. Furthermore, cognition is essential for effective human communication, including processing information given by others and responding appropriately to others. Cognitive impairment and memory loss, which tend to decline with age [37], are characteristics of diseases such as dementia. Dementia is a clinical syndrome characterized by a loss of cognitive functioning (thinking, remembering, and reasoning) that interferes with a person's daily life and activities [25]. In 2015, it was estimated that nearly 47 million people worldwide were affected by dementia, and the number continues to rise, which is expected to reach 75 million by 2030 and 131 million by 2050 [51]. Alzheimer's disease (AD) is the most common form of dementia, accounting for around 50% -70% of cases. In other words, cognitive decline is the most common cause of AD in older adults. A group analyzed the prevalence and incidence of AD in European countries by using meta-analyses. They found the prevalence of AD in Europe was estimated at 5.05% and the incidence was 11.08 per 1000 persons per year [38]. For AD or other dementia, there is no cure or treatment that substantially relieves symptoms. AD can profoundly impact the lives of patients and their families over many years, but it doesn't stop there. AD can also have a significant economic impact on individuals and their families. According to statistics, the cumulative cost of an AD patient starting from primary care to specialist care in Sweden has exceeded €5,000 [51]. However, this cost is only a small fraction of an AD patient's lifetime medical and nursing costs. The total social cost of dementia in Europe in 2010 was estimated to be between US \$238.6 billion and €105.6 billion [53]. And by 2050, the overall mean cost burden for AD and Parkinson's Disease (PD) is estimated to be €357 billion [34]. The global economic cost of dementia was estimated at more than US \$600 billion in 2010 [50], in excess of US \$1 trillion in 2020 and forecast to double by 2030 [55].

Declining cognitive functions has become one of the primary concerns for European countries. Explosive growth in care costs and associated societal burdens on AD, PD, and other dementia challenges European countries' health care and long-term care systems. Since the primary risk factor for cognitive decline is age, as life expectancy increases worldwide, the prevalence of the disease is increasing dramatically. However, Cognitive impairment has multifactorial etiology, and partial risk factors are modifiable [51]. In 2011, Barnes D E and Yaffe K identified the evidence for seven potentially modifiable risk factors associated with AD: diabetes, midlife hypertension, midlife obesity, smoking, depression, cognitive inactivity or low educational attainment, and physical inactivity. They also predicted that if all seven risk factors were reduced by 10 %-25 %, as many as 1.1 to 3 million cases of AD could be prevented globally, and 184,000 to 492,000 cases in the United States [3]. In 2014, [39] suggested that around one-third of AD cases might be attributed to potentially modifiable risk factors.

This thesis examines the association between age and cognitive decline among people, especially the elderly, living in different European countries. We model the cognitive function trajectory as a function of age to see the overall pattern and if there is a pattern change. In addition, we intend to identify all aspects of risk factors from the overall design and the pattern change. To this end, we used the data from the Survey of Health, Aging, and Retirement in Europe. SHARE is a research infrastructure providing internationally comparable longitudinal data on various health, economic, and social factors [13].

Longitudinal studies employ continuous or repeated methods for tracking the same individual to detect any changes that might occur over a long period, usually, years or decades [19]. Longitudinal data are collected naturally on any combination of exposures and outcomes, without any external influence [18]. The type of study is most commonly used in medicine, economics, medical sciences and epidemiology. It's useful for assessing risk factors, disease progression, and treatment outcomes over different lengths of time. Mixed-effects regression (MER) is one of the preferred methods for longitudinal data, which allows time-invariant(s.g., gender) and time-varying predictors (s.g., age) [36]. Furthermore, MER explicitly models individual changes over time and is flexible in repeated measures without requiring the same number of observations for each subject. Nevertheless, MER is not appropriate for our study, we instead propose an estimation procedure based on principal component bases and extend FPCA to the longitudinal setting. In MER models, we assume that our outcome variable is a linear function of age, and then there will only be a positive or negative linear relationship between them. We believe that age-based cognitive function trajectories should be more complex with a nonlinear quadratic functional format.

Functional data analysis (FDA) is a statistical field that studies models and analysis methods for data recorded over continuous time for each subject. Likewise, it can be described as the study of a sample of trajectories or time courses. In FDA, where we consider repeated observations over time courses and do not rely on any stationarity assumptions, FDA is particularly well suited for analyzing temporal dynamics and longitudinal data that are abundantly found in applications such as biomedicine. In general, any time-related data that is repeatedly observed across independent individuals or units can be analyzed using FDA's

methods. A basic paradigm of the FDA is that the observed data is regarded as an independent and identically distributed random sample from a stochastic process. However, it is also possible to incorporate dependencies between the realizations of stochastic processes. The underlying stochastic process that is assumed to lie in L_2 space is generally assumed to be smooth over a continuum, which is the target interest of FDA [23] [28]. Functional principle component analysis is one of the most popular multivariate analysis techniques for the extraction of information from the FDA.

FPCA is a powerful tool used to model longitudinal data observed at different time points. FPCA is based on the principle of component analysis(PCA); however, rather than using variables, the FPCA uses functions. It is a dimension reduction method to decompose the latent stochastic process into a linear combination of FPCs [54]. And it represents functional data in the most parsimonious way, which maximizes the variation in the randomly observed curves. The top few FPCs explain most of the variability in the underlying stochastic process [44]. FPCA can work with sparse data and does not require that the observations are taken at the same time points. In this regard, FPCA offers significant advantages for a better understanding of trends. FPCA has successfully been applied to real life, such as weather and climate prediction, fetal movement monitoring data [52], genetic growth [27] and child growth study [30].

In this thesis, because the number of repeated measurements for cognitive function per elderly is irregular and small, FPCA is performed by a nonparametric method to analyze and characterize cognitive trajectory data. Unlike classic FPCA, which requires many regularly spaced measurements to express the random curves per subject, this method works well with a small number of repeatable measurements per subject. We assume the repeatable measurements are randomly located with a random number of repetitions for each subject and are determined by an underlying smooth random trajectory plus measurement errors. Even if only one or few measurements are available for a subject, it can estimate individual smooth trajectories satisfactorily. This approach assumes the repeated measures are determined by an underlying smooth random (subject-specific) trajectory plus measurement errors [45] and principal component scores (FPCs) [54]. Simultaneously, like classic FPCA, it is suitable for extracting the pattern of the entire cognitive abilities as a function that will otherwise be lost with applying some traditional statistical techniques. By treating the whole curve as a single entity, there is no concern about correlations between repeated measurements. We assumed that an underlying functional relationship governs the data. In addition to the computational advantages, we are thus able to extract the main differences between subjects in their average cognitive function and how their cognitive function over time. We also detect the associations between one or more factors and longitudinal cognitive growth data.

The primary aim of this study is to characterize individual cognitive trajectories of the elderly using FPCA and identify the risk factors by using machine learning models (including elastic-net, decision tree and random forest). The results of FPCA show that cognitive function is an up-and-down process, divided into two stages (early and late decline), with an increase in cognitive function at about age 70, and then an accelerated decline after age 75. Based on the results of these machine learning models, we found many

preventable risk factors, such as alcohol drinking behaviour, physical activity, household income level, and education level.

The remainder of the thesis is organized as follows. Section 1 is the introduction. In Section 2, we present the data, including the source of data, prime outcome, covariates and study sample. The FPCA, Cross-Validation, Elastic-net, Decision tree, and Random Forest methods are described in Section 3. The result of FPCA and risk factors are presented in Section 4. Section 5 is the conclusion.

2 The data

2.1 Source of data

Our sample data is from SHARE. The baseline study of SHARE took place in 2004 and continued today. With the exception of the SHARE, Corona Survey collected via telephone in 2020, all other SHARE data collection consisted of face-to-face computer-assisted interviews. 140,000 people aged 50 or older from various regions of Europe and Israel participated and conducted 530,000 in-depth interviews. The baseline wave involved 12 countries (Austria, Belgium, Denmark, France, Germany, Greece, Israel, Italy, Netherlands, Spain, Sweden, and Switzerland). Switzerland and Belgium contribute the least and the most respondents at around 3.33% and 12.68% of the whole population, respectively. Approximately 8% to 10% of respondents are from other countries. In contrast, Spain and Italy have seriously ageing populations, with the average respondent reaching over 65 years old, while the average age of respondents in Greece is relatively low, around 62 years old. Respondents from other countries are about 63 to 64 years old. Overall, the male-to-female ratio in these participating countries is around 5 to 5. Several other countries have been added in the following waves. Until now, 28 European countries and Israel joined SHARE. The persons aged 50 years and over are the target population of SHARE. New participants will be enrolled in each wave as refreshment samples to compensate for the dropout of participants. The collection method in wave 3 differs from other waves; it is called SHARELIFE. The questioning method in SHARELIFE is based on the Life History Calendar (LHC). Respondents' lives are graphically represented by a grid that automatically populates during the interview. The LHC helps respondents remember primarily by asking about life events that are likely to be reflected accurately. Usually includes the respondent's child's name, date of birth and partner's history. Since our outcome variable has not been provided in the third wave, the sample data include wave 1 (2011), wave2 (2011), wave4 (2013), wave5 (2015), wave6 (2017), wave7 (2019) and wave8 (2020). Wave 1 was selected as the baseline, and wave 8 was a new wave when collection took place in 2020.

2.2 Primary Outcome

Cognition function was assessed using the "ten words list learning" test in SHARE. The test is conducted with immediate recall and delayed recall. Immediate recall was the number of recalled words after the interviewer read a list of 10 words, so the score of immediate recall ranges from 0 to 10. At the end of the test, the interviewer will ask the participants to recall the words again from the list as the delayed recall

score(DRS). From Table 1, the mean of DRS decreases as the wave increases, from 3.44 in wave 1 to 2.54 in wave 8, showing a downward trend; as we mentioned in the introduction, cognitive function is closely related to age. We will use the normalized delayed recall score as our outcome variable and only keep the subjects without any missing value in DRS. According to our introduction to cognitive function in section 1, DRS belongs to cognitive mechanics, which is sensitive to age differences.

2.3 Covariates

SHARE is an extensive database that covers thousands of variables; we initially identified covariates from multi-aspect, including demographic information, household composition, social support, network, childhood conditions, health, behavioural risk, work, and money. Our descriptive statistics for the final covariates across waves are depicted in Tables 2.1 and 2.2. Note there are 944 people in each wave, no missing values in the wave1, but other waves might have it, and the number in the bracket is the percentage of an index for a variable (Table 2.1). The number in the bracket after the description is the index in the dataset. For example, in the original dataset, Male(1) means Male represent by index 1.

In descriptive Table 2.1, we found that the participants were more female (58%), and the elderly living alone were increasing, from 29% in wave 1 to 54% in wave 8. At the same time, the health of the elderly is declining, the number of hospitalisations is increasing, and some activities cannot be completed, which is in line with the laws of nature. It is worth noting that when self-rated writing, reading and orientation time skills have dropped significantly, the numeracy score of the elderly has not shown a significant drop; it remained unchanged. From the descriptive Table 2.2, we found that the cognitive function of the elderly decreased with age, which was consistent with expectations. In addition, BMI, maximum grip strength, and income levels all decreased over time.

Table 2.1: Descriptive Statistics Across Waves (discrete variable)

Variable	wave 1	wave 2	wave 4	wave 5	wave 6	wave 7	wave 8
gender							
Male(1)	401(42)	401(42)	401(42)	401(42)	401(42)	401(42)	401(42)
Female(2)	543(58)	543(58)	543(58))	543(58)	543(58)	543(58)	543(58)
partnerinhh							
Living with spouse/partner(1)	670(71)	642(68)	581(62)	545(58)	506(54)	475(50)	432(46)
Living w.o spouse/partner(3)	274(29)	302(32)	363(38)	399(42)	438(46)	469(50)	512(54)
Marital status(mstat_m)							
married and living together with spouse(1)	645(68)	25(45)	16(27)	24(34)	24(37)	10(17)	7(11)
registered partnership(2)	10(1)	0(0)	0(0)	0(0)	0(0)	0(0)	0 (0)
married, living separated from spouse(3)	7(1)	0(0)	0(0)	0(0)	1(2)	1(2)	1(2)
never married(4)	38(4)	1(2)	1(2)	1(1)	2(3)	0(0)	1(2)
divorced(5)	54(6)	0(0)	3(5)	2(3)	1(2)	1(2)	2(3)
widowed(6)	190(20)	29(53)	40(67)	44(62)	37(57)	46(79)	50(82)
Current job situation(ep005)							
employed(1)	39(4)	18(2)	13(1)	7(1)	6(1)	8(1)	4(0)
unemployed,retired(0)	904(96)	918(98)	914(99)	919(99)	913(99)	902(99)	880(100)
Number of Child(child)							

0(1)	76(11)	79(11)	81(11)	87(12)	90(12)	88(12)	88(11)
1 to 3(2)	511(72)	519(72)	535(72)	530(72)	557(73)	540(73)	568(73)
over 4(3)	126(18)	120(17)	124(17)	117(16)	120(16)	113(15)	121(16)
Number of Received help(ghelp)							
0(0)	758(80)	789(84)	758(80)	741(78)	685(73)	584(65)	557(59))
1 to 3(1)	186(20)	155(16)	186(20)	203(22)	259(27)	318(35)	387(41)
Number of Given help(ghelp)							
0(1)	609(65)	645(68)	553(75)	558(76)	750(79)	735(81)	808(86)
1 to 3(2)	335(35)	299(32)	187(25)	179(24)	194(21)	167(19)	136(14)
Self-perceived health(health)							
Excellent(1)	112(12)	94(10)	66(7)	56(6)	51(5)	39(4)	34(4)
Good(2)	636(67)	567(60)	581(62)	547(58)	526(56)	505(53)	445(47)
Poor(3)	196(21)	283(30)	297(31)	341(36)	367(39)	400(42)	465(49)
Hospital stay(hospital_m)							
No(0)	849(90)	814(86)	791(84)	780(84)	775(82)	757(80)	724(77)
Yes(1)	95(10)	130(14)	151(16)	164(17)	169(18)	187(20)	220(23)
Drinking behavior(drink)							
not at all(1)	273(31)	249(29)	265(31)	281(33)	NaN	NaN	NaN
less than twice a month(2)	20(2)	72(8)	73(9)	78(9)	NaN	NaN	NaN
less than four days a week(3)	293(34)	246(29)	228(27)	225(27)	NaN	NaN	NaN
almost every day(4)	281(32)	286(34)	277(33)	262(31)	NaN	NaN	NaN
ever smoking(esmoked_m)							
No(0)	589(62)	592(63)	592(63)	593(63)	593(63)	567(63)	607(64)
Yes(1)	355(38)	351(37)	351(37)	351(37)	351(37)	335(37)	337(36)
N of chronic diseases(chronic)							
0(1)	198(21)	191(20)	159(17)	145(15)	134(14)	117(12)	100(11)
1(2)	316(33)	299(32)	289(31)	256(27)	235(25)	231(24)	226(24)
over 1(2)	430(46)	454(48)	496(53)	543(58)	575(61)	596(63)	618(65)
N of doctor visit in 12 month(doctorS)							
0 to 3(1)	430(46)	393(42)	343(36)	331(35)	328(35)	327(35)	290(31)
4 to 12(2)	417(44)	426(45)	474(50)	468(50)	497(53)	485(51)	492(52)
13 to 98(3)	97(10)	125(13)	127(13)	145(15)	119(13)	132(14)	162(17)
depression scale(depression)							
0 to 3(1)	691(76)	730(79)	671(74)	633(70)	644(71)	590(68)	586(64)
4 to 12(2)	220(24)	190(21)	237(26)	275(30)	265(29)	273(32)	330(36)
Mobility Index(mob)							
0(1)	473(50)	496(53)	420(45)	378(40)	375(40)	326(35)	266(28)
1 to 4(2)	471(50)	448(47)	522(55)	566(60)	569(60)	618(65)	678(72)
Eyesight reading(reading)							
poor(0)	180(19)	161(17)	174(18)	214(23)	182(19)	161(18)	208(22)
good(1)	764(81)	783(83)	770(82)	730(77)	762(81)	741(82)	736(78)
Self-rated writing skills(writing)							
poor(0)	169(18)	167(18)	167(18)	166(18)	166(18)	219(24)	222(24)
good(1)	775(82)	777(82)	777(82)	778(82)	778(82)	683(76)	722(76)
Self-rated reading skills(reading)							
poor(0)	143(15)	144(15)	144(15)	143(15)	143(15)	153(17)	155(16)
good(1)	801(85)	800(85)	800(85)	801(85)	801(85)	749(83)	789(84)
hearing							
poor (0)	171(18)	188(20)	215(23)	261(28)	260(28)	271(30)	306(32)
good(1)	773(82)	756(80)	729(77)	683(72)	684(72)	631(70)	638(68)
Large Muscle Index(muscle)							
0(1)	567(60)	597(63)	522(55)	487(52)	490(52)	432(46)	64(36)
1 to 2(2)	229(24)	203(22)	225(24)	245(26)	219(23)	246(26)	34(19)
3 to 4(3)	148(16)	144(15)	197(21)	212(22)	235(25)	266(28)	80(45)

Vigorous activities(vigorous)							
never(1)	348(37)	395(42)	463(49)	508(54)	529(56)	535(59)	124(70)
more than once(2)	596(63)	548(58)	480(51)	436(46)	415(44)	367(41)	54(30)
active of Daily Living Index(active)							
0(1)	899(95)	886(94)	857(91)	846(90)	820(87)	809(86)	714(76)
1 to 5(2)	45(5)	58(6)	85(9)	98(10)	124(13)	135(14)	230(24)
Instrumental Activities Index(active2)							
0(1)	844(89)	848(90)	794(84)	762(81)	700(74)	670(71)	536(57))
1 to 5(2)	100(11)	96(10)	148(16)	182(19)	244(26)	274(29)	408(43)
Fine Motor Skills Index(finemotor)							
1	891(94)	882(93)	852(90)	837(89)	829(88)	820(87)	135(76)
2	50(5)	55(6)	85(9)	96(10)	94(10)	102(11)	30(17)
3	3(0)	7(1)	7(1)	10(1)	18(2)	18(2)	11(6)
4	0(0)	0(0)	0(0)	0(0)	0(0)	0(0)	0(0)
Gross Motor Skills Index(grossmotor)							
1	851(90)	843(89)	790(84)	766(81)	747(79)	674(71)	96(54)
2	74(8)	70(7)	107(11)	107(11)	100(11)	158(17)	39(22)
3	15(2)	23(2)	37(4)	52(6)	66(7)	70(7)	22(12)
4	2(0)	7(1)	8(1)	10(1)	22(2)	31(3)	13(7)
5	2(0)	1(0)	2(0)	9(1)	8(1)	10(1)	8(4)
Score of orientation in time test(orient)							
1	0(0)	0(0)	0(0)	0(0)	4(0)	9(1)	30(3)
2	0(0)	4(0)	5(0)	5(0)	6(1)	13(1)	19 (2)
3	16(2)	6(1)	8(0)	10(1)	15(2)	31(3)	46(5)
4	108(11)	98(10)	100(11)	102(11)	111(12)	125(14)	163(17)
5	820(87)	836(89)	831(88)	827(88)	808(86)	724(80)	686(73)
Numeracy Score(percentage)(numeracy_m)							
1	44(5)	45(5)	46(5)	46(5)	46(5)	43(5)	46(5)
2	154(16)	156(17)	156(17)	155(16)	155(16)	151(17)	155(16)
3	312(33)	274(29)	273(29)	273(29)	272(29)	259(29)	272(29)
4	290(31)	303(32)	302(32)	303(32)	304(32)	286(32)	304(32)
5	144(15)	166(18)	167(18)	167(18)	167(18)	163(18)	167(18)
mother alive							
Yes(1)	70(7)	45(23)	15(2)	15(2)	9(47)	4(0)	0(0)
No(5)	867(93)	153(77)	750(98)	929(98)	10(53)	927(100)	20(100)
father alive							
Yes(1)	9(1)	3(2)	1(0)	1(0)	1(12)	0(0)	0(0)
No(5)	930(99)	146(98)	760(100)	939(100)	7(88)	934(100)	14(100)
country							
Austria	38(4)	38(4)	38(4)	38(4)	38(4)	38(4)	38(4)
Belgium	137(15)	137(15)	137(15)	137(15)	137(15)	137(15)	137(15)
Denmark	88(9)	88(9)	88(9)	88(9)	88(9)	88(9)	88(9)
France	126(13)	126(13)	126(13)	126(13)	126(13)	126(13)	126(13)
Germany	95(10)	95(10)	95(10)	95(10)	95(10)	95(10)	95(10)
Italy	133(14)	133(14)	133(14)	133(14)	133(14)	133(14)	133(14)
Sweden	137(15)	137(15)	137(15)	137(15)	137(15)	137(15)	137(15)
Switzerland	75(8)	75 (8)	75(8)	75(8)	75(8)	75(8)	75(8)
Spain	115(12)	115(12)	115(12)	115(12)	115(12)	115(12)	115(12)

Table 2.2: Descriptive Statistics Across Waves (continues variable)

Variable	Summary	wave 1	wave 2	wave 4	wave 5	wave 6	wave 7	wave 8
Immediate recall Score	mean(sd)	4.85(1.68)	4.92(1.67)	4.93(1.58)	4.83(1.63)	4.74(1.65)	4.54(1.69)	4.03(1.77)
Age for time interview	mean(sd)	69.31(3.74)	71.54(3.78)	75.88(3.74)	77.83(3.76)	79.78(3.76)	81.86(3.75)	84.48(3.76)
BMI	mean(sd)	26.41(3.89)	26.50(3.99)	26.45(4.08)	26.60(4.15)	26.49(4.30)	26.23(4.29)	25.80(4.43)
Income(thousand)(thinc_m)	median	25.11	21.78	24.24	23.58	21.64	20.42	21.00
	IQR	26.67	22.95	25.46	22.74	20.04	19.45	18.87
Education year	mean(sd)	9.38(4.81)	10.38(4.90)	10.37(4.92)	10.38(4.89)	10.38(4.87)	10.37(4.88)	10.37(4.88)
Maximum of grip strength(maxgrip)	median	31	30	29	27	27	26	24
	IQR	16	16	15	15	14	14	13
Household size(hhsize)	median	2	2	2	2	2	2	2
	IQR	1	1	1	1	1	1	1
Number of siblings alive	median	1	1	1	1	1	1	2
	IQR	2	2	2	2	2	1	3

2.4 Study sample

Due to filters/routing or an abandoned interview, there are plenty of missing values in the main SHARE database. Besides, apart from the above reasons, SHARE uses missing codes to represent the missing values, such as -1 , -2 , and -99 . We regard all these missing codes as missing values. Independently of the chosen imputation method. SHARE provide five multiple imputations [42] of the missing values on some variables by using the simple hot-deck method or jointly by the fully conditional specification method (FCS) [49] for users to account for additional variability caused by the imputation process when evaluating the accuracy of their estimators. Hot-deck imputations are carried out separately by country, while FCS imputations are carried out by country and sample type. In our study, we use some variables that have been imputed by SHARE to reduce the proportion of missing values. For continuous variables, we calculate the mean of the five multiple imputations. And for discrete variables, we use the most frequent value.

Researchers in different fields often wonder what percentage of missing data should be removed. Yet, there is generally no established cut-off value for an acceptable portion of missing data in a dataset. By referring to various literature, we find that Schafer asserts that a missing rate of 5% or less is irrelevant [43] and Bennett believes that when more than 10% of the data are missing, the statistical analysis is likely to be biased [5]. Considering our dataset, we think 10% is an appropriate threshold. We further filter out variables if there is more than 10% missing data [33].

We suspected an accelerated decline in cognitive function after the age of 65, So we kept only participants whose baseline age is from 65 to 80. Among these, we find the sample size for the age over 95 is too small, which causes the outliers in the FPCA model. Therefore we remove the age over 95 in wave 8. In addition, we filter out subjects with DRS less than zero and do not consider interviewees with missing values in DRS. Since many respondents did not participate in all wave interviews, we exclude those respondents who did not have all 7 measurements (all waves except wave 3). The FPCA method sample includes 944 individuals with 6608 measurements. After filtering out variables with missing data, the sample size reduces to 713.

3 Methods

3.1 Sparse Functional principal component analysis

In sparse FPCA, sparse functional data as noisy sampled points are assumed to be independent realizations of a smooth random function with the mean is $E(x_i(t)) = \mu_t$ and $\text{cov}(x_i(s), x_i(t)) = G(s, t) = \sum_k \lambda_k \phi_k(s) \phi_k(t)$, $t, s \in \mathbf{T}$, where λ_k are non-increasing eigenvalues. The t and s are index variables as time which belong to the closed time interval \mathbf{T} . In classical FPCA, $\xi_{ik} = \int (x_i(t) - \mu(t)) \phi_k(t) dt$ and the i^{th} random curve is expressed as $X_i(t) = \mu(t) + \sum_k \xi_{ik} \phi_k(t)$. $t \in \mathbf{T}$, where ξ_{ik} is the k^{th} FPCs for i^{th} random curve with $E(\xi_{ik}) = 0$, $\text{Var}(\xi_{ik}) = \lambda_k$, $\sum_k \lambda_k < \infty$, $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$, and $\phi_k(t)$ is the k^{th} eignfunction. X_i can be represented by the sequence of $\xi_{i1}, \xi_{i2}, \dots, \xi_{ik}$. For any k , the first k term is the best k -dimensional linear approximation for $X(t)$ in L_2 space, which explains the highest proportion of variance in data with a given number of components.

Now we consider adding uncorrelated measurement errors with mean zero and constant variance σ^2 to reflect additive measurement errors and the model will become [54]:

$$Y_{ij} = X_i(T_{ij}) + \epsilon = \mu(T_{ij}) + \sum_{k=1}^{\infty} \xi_{ik} \phi_k(T_{ij}) + \epsilon_{ij} \quad (3.1)$$

where $T_{ij} \in \mathbf{T}$, ϵ_{ij} is the i.i.d. measurement error and is also independent of ξ_{ik} , where $i = 1, \dots, n$, $j = 1, \dots, N_i$, $k = 1, 2, \dots$. N_i is the number of measurements on the i^{th} subject.

Note Y_{ij} is the j^{th} observation of the random function $X_i(t)$, made at a random time T_{ij} . In our study, Y_{ij} is a j^{th} cognitive data point for i^{th} person and T_{ij} is the age in the j^{th} measurement for i^{th} person, where i and j are integers. $j \in \{1, 2, 3, 4, 5, 6, 7\}$ and $i \in \{1, 2, 3, 4, \dots, 944\}$

In sparse FPCA, we use local linear smoothers for function and surface estimation. The local lines and planes are fitted by weighted least squares. Mean, covariance and eigenfunctions are all assumed to be smooth [22]. The mean function $\hat{\mu}$ is estimated based on the pooled data from all individuals. The local linear scatterplot smoother for $\hat{\mu}$ is minimizing

$$\sum_{i=1}^n \sum_{j=1}^{N_i} \mathbf{k}_1\left(\frac{T_{ij} - t}{h_\mu}\right) (Y_{ij} - \beta_0 - \beta_1(t - T_{ij}))^2 \quad (3.2)$$

with respect to β_0, β_1 . The estimate of $\mu(t)$ is $\hat{\mu}(t) = \hat{\beta}_0(t)$, \mathbf{k}_1 is a kernel function of order (\mathbf{v}, l) and is compactly supported, $\|\mathbf{k}_1\|^2 = \int \mathbf{k}_1^2(u) du < \infty$, where \mathbf{v} is a multi-index $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2)$. $\mathbf{v}_1, \mathbf{v}_2, l$ are the given integers, with $0 \leq \mathbf{v}_1 + \mathbf{v}_2 \leq l$

The local linear surface smoother for $G(s, t)$ is defined by minimizing

$$\sum_{i=1}^n \sum_{1 \leq j \neq l \leq N_i} \mathbf{k}_2\left(\frac{T_{ij}-s}{h_G}, \frac{T_{il}-t}{h_G}\right) (G_i(T_{ij}, T_{il}) - f(\beta, (s, t), (T_{ij}, T_{il})))^2 \quad (3.3)$$

where $f(\beta, (s, t), (T_{ij}, T_{il})) = \beta_0 + \beta_{11}(s - T_{ij}) + \beta_{12}(t - T_{il})$. Minimization is with regard to $\beta = (\beta_0, \beta_{11}, \beta_{12})$. \mathbf{k}_2 is a kernel function of order (\mathbf{v}, l) and is compactly supported, $\|\mathbf{k}_2\|^2 = \int \int \mathbf{k}_2^2(u, \mathbf{v}) du d\mathbf{v} < \infty$. Usually, the smoothing parameter for this surface smoothing step is chosen by one-curve-leave-out cross-validation.

The covariance function between Y_{ij} and Y_{il} is given by

$$\text{cov}(Y_{ij}, Y_{il} | T_{ij}, T_{il}) = \text{cov}(X_i(T_{ij}), X_i(T_{il})) + \sigma^2 \delta_{jl} \quad (3.4)$$

where $\delta_{jl} = 1$ if $j = l$ and 0 otherwise.

Now let $G_i(T_{ij}, T_{il}) = (Y_{ij} - \hat{\mu}(T_{ij}))(Y_{il} - \hat{\mu}(T_{il}))$ be the raw covariances, where $\hat{\mu}(t)$ is the estimated mean function. Then the $E[G_i(T_{ij}, T_{il})]$ becomes

$$E[G_i(T_{ij}, T_{il})] \approx \text{cov}(X_i(T_{ij}), X_i(T_{il})) + \sigma^2 \delta_{jl} \quad (3.5)$$

By estimating the diagonal covariance function $V(t) = G(t, t)$ with the local linear smoother on the diagonal raw covariances $G_i(T_{ij}, T_{ij})$ (obtained from equation (3.2) by using $\{G_i(T_{ij}, T_{ij})\}$ as input) [29] [21] [24], if $\sigma^2 > 0$, the σ^2 of the measurement errors in Equation (3.1) is estimated by

$$\hat{\sigma}^2 = \frac{2}{|\mathbf{T}|} \int_{\mathbf{T}_1} [\hat{V}(t) - \hat{G}(t, t)] dt \quad (3.6)$$

where $|\mathbf{T}|$ denote the length of \mathbf{T} and \mathbf{T}_1 is the interval $\mathbf{T}_1 = [\inf\{x : x \in \mathbf{T}\} + |\mathbf{T}|/4, \sup\{x : x \in \mathbf{T}\} - |\mathbf{T}|/4]$

The eigenfunctions and eigenvalues correspond to the $\hat{\phi}_k$ and $\hat{\lambda}_k$ can be estimated by the following equation with the constraints $\|\hat{\phi}_k(t)\| = 1$ and $\langle \hat{\phi}_k(t), \hat{\phi}_m(t) \rangle = 0$, for $m < k$. $\hat{G}(s, t)$ denote the smooth surface estimate of $G(s, t) = \text{cov}(X(s), X(t))$ (see Equation (3.3))

$$\int_{\mathbf{T}} \hat{G}(s, t) \hat{\phi}_k(s) ds = \hat{\lambda}_k \hat{\phi}_k(t) \quad (3.7)$$

For the classical FPCA, FPC scores are estimated by $\xi_{ik} = \int (x_i(t) - \mu(t)) \phi_k(t) dt$, which works well for the sufficient density of the grid of measurements for each subject. However, for sparse FPCA, since the Y_{ij} are only available at discrete random times T_{ij} , this numerical integration can not provide a reasonable approximation. The FPCs for sparse FPCA should be estimated under the assumption that ξ_{ik} and ϵ_{ij} are jointly Gaussian, and the best prediction is given by

$$\tilde{\xi}_{ik} = E(\xi_{ik} | Y_i) = \lambda_k \phi_{ik} \sum_{Y_i}^{-1} (Y_i - \mu_i) \quad (3.8)$$

where $Y_i = (Y_{i1}, \dots, Y_{iN_i})$, $\sum_{Y_i} = \text{cov}(Y_i, Y_i)$, $\phi_{ik} = (\phi_k(T_{i1}), \dots, \phi_k(T_{iN_i}))$ and $\mu_i = (\mu(T_{i1}), \dots, \mu(T_{iN_i}))$, Note the true values of λ_k , ϕ_{ik} , \sum_{Y_i} and μ_i are unknown. $\tilde{\xi}_{ik}$ is the best linear prediction of ξ_{ik} , given the information from the i^{th} subject, no matter whether the Gaussian assumption holds or not.

By substituting estimates of μ_i , λ_k , ϕ_{ik} and \sum_{Y_i} from the entire data ensemble, Yao, Muller and Wang (2005) proposed to use the conditional expectation to estimate scores

$$\hat{\xi}_{ik} = \hat{\lambda}_k \hat{\phi}_{ik} \sum_{Y_i}^{-1} (Y_i - \hat{\mu}_i) \quad (3.9)$$

The prediction of trajectory $X_i(t)$ is given by

$$\hat{X}_i^K(t) = \hat{\mu}(t) + \sum_{k=1}^K \hat{\xi}_{ik} \hat{\phi}_k(t) \quad (3.10)$$

The K is chosen by cross-validation based on the one-cure-leave-out prediction error [54]. This conditioning method provides the best predictors under Gaussian assumptions and works in the presence of both measurement errors and sparsity.

All statistical analyses were performed using R studio version 4.0.4 and FPCA was performed using the *fdapace* package [21] [54] [1] [20]. The first step toward using the FPCA function in the *fdapace* package is to restructure the data so that each respondent's time and DRS data are stored as lists in separate columns, where each row contains all of the data for a single id. The work-flow for sparse FPCA is as follows:

- Calculate the smoothed mean $\hat{\mu}$ (local linear smoothing) by aggregating all available cognitive curves together (see Equation (3.2)).
- Calculate each cognitive curve's own raw covariance separately and aggregate all these raw covariances to form the sample raw covariance (see Equation (3.4) and (3.5)).
- Use the off-diagonal raw covariances to estimate the smooth covariance (see Equation (3.3)).
- Obtain the eigenfunctions $\hat{\phi}$ and eigenvalues $\hat{\lambda}$ by doing eigenanalysis on the smooth covariance and then project the smooth covariance on a positive semi-definite surface (see Equation (3.7)) [26].
- Use conditional expectation to estimate the $\hat{\xi}$ (see Equation (3.8) and (3.9)).

In functional data analysis, FPCA plays a vital role. The top FPCs explain major sources of total variation among the whole dataset. We use FPCA to catch major variance, avoid overfitting, and smooth the fitted curves [44]. The FPCA method can accurately capture the time fluctuation characteristics of cognitive function, especially the changing direction and form in time, which provides a scientific basis for modelling the cognitive function of the elderly. In the next section, we introduce some machine learning models and methods, namely cross-validation, Lasso, ridge, elastic network model, decision tree and random forest model.

3.2 Statistical learning

3.2.1 Shrinkage Models

After applying the FPCA method for dimension reduction, we use shrinkage models to identify the risk factors. In the machine learning field, it is called feature selection. Our dependent variable is FPC1 or

FPC2 (y_i) which we obtain from the FPCA technique and independent variables (x_{ij}) are the covariates in Table 1 and 2. Note that since after removing all missing values there are 713 individuals and 39 covariates remaining, so $i \in \{1, 2, 3, \dots, 713\}$ and $j \in \{1, 2, 3, \dots, 39\}$. After fitting the shrinkage model with significant coefficients(β_j), the remaining covariates are the risk factors we want to recognise.

Ridge and Lasso Regression

The ridge regression and Lasso regression are the two best-known techniques for the shrinkage of the regression coefficients. By adding the L_1 or L_2 penalty into the loss function, both methods could shrink the regression coefficients toward zero. These two methods are very similar to least squares ($RSS = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$), except that the coefficients are estimated by minimizing a slightly different quantity. The ridge regression provides the minimum lambda to minimize the

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (3.11)$$

where λ is a tuning parameter that greater or equal to zero. The shrinkage penalty $\lambda \sum_{j=1}^p \beta_j^2$ has the effect of shrinking the estimates of β_j towards zero, but it never shrinks the β_j exactly to zero. The parameter λ is used to control these two terms' relative influence on the regression coefficients' estimation. That is, when λ is zero, the penalty term has no effect, and when λ approaches infinite, the penalty becomes more impactive, ridge regression coefficient tends to be zero. Ridge regression has the obvious disadvantage that the final ridge regression model will contain all predictors. However, the Lasso regression can overcome this disadvantage.

The Lasso regression provides the minimum lambda to minimize the

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (3.12)$$

L_1 penalty equal to $\lambda \sum_{j=1}^p |\beta_j|$ can estimate all the coefficients exactly to zero when the tuning parameter λ is sufficiently large. Therefore, the Lasso performs variable selection, and the Lasso model is much easier to interpret than the ridge. The λ usually is selected by cross-validation methods. By choosing a grid of λ values, we select the λ for which the cross-validation error is the smallest. Lasso also has obvious shortcomings. When there is a multicollinearity problem, Lasso randomly selects a multicollinear variable, so Lasso cannot overcome the multicollinearity problem. The shrinkage regressions (e.g. ridge, Lasso regression) usually significantly reduce the prediction variance and achieve the purpose of coefficient shrinkage and variable selection, but they all have limitations.

Elastic Net

The elastic net is a hybrid between the Lasso and ridge model for improving the shortcomings of these two techniques; it uses the penalties as a combination of a L_1 norm penalty and a L_2 norm penalty. Ridge

with L_2 norm $= \sum_{j=1}^p \beta_j^2$ includes all the predictors in the final model and will not perform feature selection. Lasso with L_1 norm $= \sum_{j=1}^p |\beta_j|$ only takes a few samples for high-dimensional data. Like ridge and Lasso, we attempt to minimize the residual sum of squares with different penalty terms:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \left[(1 - \alpha) \left(\sum_{j=1}^p \beta_j^2 \right) / 2 + \alpha \sum_{j=1}^p |\beta_j| \right] \quad (3.13)$$

note that dividing the ridge penalty by 2 is convenient for optimization.

The new penalty is $\lambda \left[(1 - \alpha) \left(\sum_{j=1}^p \beta_j^2 \right) / 2 + \alpha \sum_{j=1}^p |\beta_j| \right]$, which balance the two penalties of Lasso and ridge and can result in a better performance on some problems. Notice that α and λ are always positive numbers. It is more useful to think of α as controlling the mixing between the two penalties and λ controlling the amount of penalization. α takes values between 0 and 1. When $\alpha = 1$, it gives Lasso and when $\alpha = 0$, it gives ridge. We use R Package *caret* to fit the elastic-net model [32]. In that package, we set up a 10 fold cross-validation strategy and use *train()* with *method* = “*glmnet*” to fit the elastic net. We allow *caret* automatically choose the best tuning parameters α and λ based on the minimizing cross-validation error.

If a group of variables is highly correlated, Lasso tends to select only one in this group and ignore the rest. In elastic-net, strongly correlated predictors tend to appear together in or removed from the model [56]. Hence, elastic-net performs feature selection and regularization simultaneously. In our case, some selected variables may be highly correlated. With the grouping effect of elastic-net regularization, those correlated variables will be selected together, and non-zero coefficients will become more interpretable.

3.2.2 Decision tree

The decision tree method is a predictive algorithm that can perform both classification and regression based on multiple covariates. Common usages of decision tree models include feature selection, assessing the relative importance of variables, handling of missing values, prediction, and data manipulation. Besides, decision trees are fundamental components of random forests. The main components of a decision tree model are nodes and branches. Each node represents a “test”, and each branch represents the result of the test. The node in the tree model shows the predicted value, the number of data points reaching this node, and the population percentage in this node.

The Figure named Decision Tree shows the prediction of the FPC1 from the categorical variables. This decision tree starts with a single node (base. country = Austria, Denmark, France, Germany, Sweden, Switzerland), which branches into possible outcomes. Each outcome leads to additional nodes, which branch off into other possibilities. 0.032 is a predicted value for the first FPCs, 713 is the number of data points reaching this node, and 100% is the percentage of the population in this node. For example, if you choose “Yes” in the first node, it comes to the node base.numeracy = 4 or 5 (with probability 58%) and if you choose “No”, it comes to the node base.writing = good (with probability 42%).

All statistical analyses about decision tree are using R Package *rpart* [48] and all visualization plots are

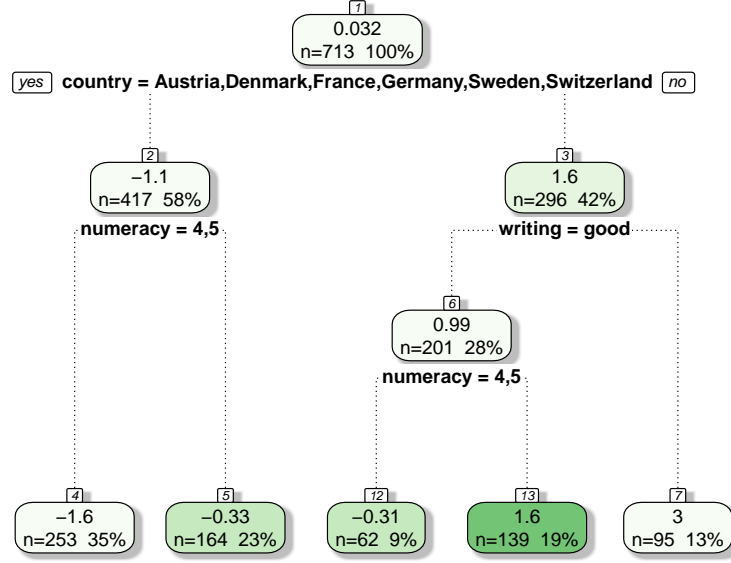


Figure 3.1: Decision Tree

using R Package *rpart.plot* [35]. Since the outcome variables, FPC1 or FPC2 are continuous. We need to build the regression tree rather than the classification tree. The process of building a regression tree is roughly two steps. First, we construct distinct and non-overlapping nodes and put each observation onto these nodes. Second, we make the same prediction for every observation that falls into the node by simply calculating the mean response values for the training observations. For example, if there are node 1 and node 2 obtained in the first step. The response mean of the training observations in the first node is 5 and in the second node is 10. Then if given observations $a_1 \in \text{node 1}$ and $a_2 \in \text{node 2}$, we will predict a_1 as 5 and a_2 as 10. We construct the node by minimizing the residual sum of square (RSS), given by

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 \quad (3.14)$$

where R_j represents j^{th} node, $y_i - \hat{y}_{R_j}$ is the mean response for the training observations within the j^{th} node. The decision tree first considers all predictors and all possible values of cutpoint for each predictor. Then it chooses the predictor and cutpoint such that the resulting tree has the lowest RSS.

The resulting decision tree might be too complex, which is likely to overfit the data, leading to poor test set performance. Therefore, we will prune a tree to obtain a subtree. Rather than considering every possible subtree, cost complexity (CP) determines the best prune way to prune the tree. We consider a sequence of

trees indexed by a nonnegative tuning parameter α . For each value of α , there corresponds a subtree $T \in T_0$ such that

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T| \quad (3.15)$$

is as small as possible. Here T_0 is the tree before pruning, and $|T|$ is the number of terminal nodes of the tree T . R_m is the subset of predictor space (rectangle) corresponding to the m^{th} terminal node. \hat{y}_{R_m} is the predicted response associated with R_m . The parameter α controls the trade-off between the complexity of the subtree and the fit to the training data. When $\alpha = 0$, $T = T_0$. Usually, we use cross-validation to select the value of α and then return it to the dataset for obtaining the subtree corresponding to α .

The final predicted value for the decision tree is FPC1 or FPC2, which obtain from the FPCA technique, and the nodes are picked from the covariates in Table 1 and 2. Our objective is to find the risk factor, which is the name of each node in the decision tree. The step for the decision tree:

- Pick the variables based on the lowest RSS and grow a large tree on the training data.
- Snipping off the least important splits based on the CP (prune the tree).
- Use cross-validation (Repeat step 1 and step 2 on all but the k^{th} fold of the training data) to choose the α that minimizes the mean squared prediction error.
- Obtain the subtree in step 2 with the α in step 3.

The decision tree model is straightforward to understand and interpret, and it can handle heavily skewed without doing transformation, and missing data without imputation [47]. However, decision trees are not stable; a tiny change in data may lead to a significant difference in the structure of the optimal decision tree. Decision trees also tend to have an over-fitting problem, which random forests can handle, but they are not as easily explained as decision trees.

3.2.3 Random forest

The Random Forest algorithm proposed by L. Breiman in 2001 is a powerful machine learning classifier that performs classification by constructing numerous decision trees by bootstrapping and aggregating their results to create the final model (called bagging). Bootstrapping is resampling the observed dataset (and of equal size to the observed dataset), each of which is obtained by random sampling with replacement from the original dataset [14] [46]. Figure 3.2 depicts how Bootstrapping works. The right part is the original dataset, and the left details are the training datasets bootstrapping. It displays a small sample size $n = 3$. Each bootstrap dataset contains 3 observations, and it chooses samples with replacements from the original dataset. $\hat{\alpha}^{*i}$ are estimated by each bootstrap dataset.

By generating B different bootstrapped training data sets, we train our method on the b^{th} bootstrapped training data set to get $\hat{f}^{*b}(x)$. After averaging all the predictions, we obtain the following equation called

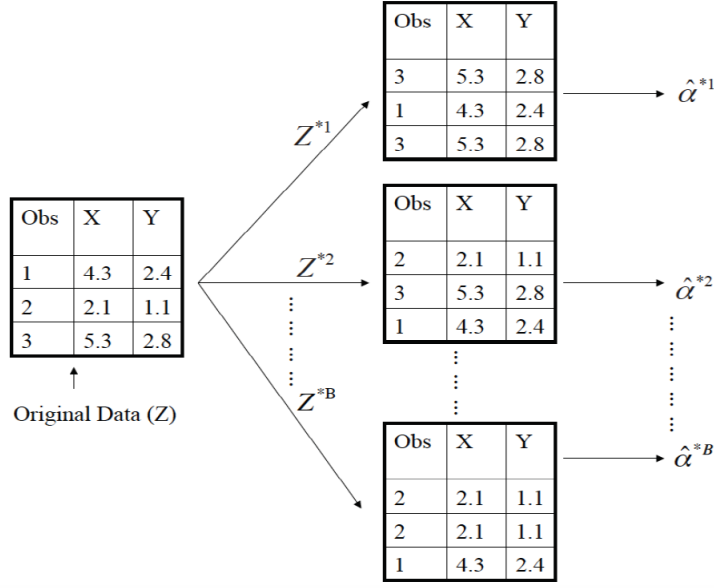


Figure 3.2: Bootstrapping

bagging. With bigger values of B , the variance of the decision trees will decrease while the computational time will grow substantially. An optimal number of trees B can be found using cross-validation to balance precision and computational cost. However, bagging contains a large number of decision trees, it is no longer possible to exhibit a process like a single decision tree, and it is no longer clear which variables are critical to the overall process. Thus, the improved predictive power of bagging comes at the expense of interpretability.

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x) \quad (3.16)$$

We build several decision trees on bootstrapped training samples in the random forest model. When building these decision trees, a split in a tree is considered at each time, and a random sample with m predictors is chosen from the full set of p predictors. Splitting allows only one of the m predictors to be used. The new sample of m predictors is taken for each split. Usually, we choose $m \approx \sqrt{p}$. Note that m is the number of predictors considered at each split, and p is the total number of predictors. The difference between bagging and the random forest is the predictor subset size m . If a random forest is built using the m equal to p , this amounts simply to bagging.

Because of the law of large numbers, random forest rectifies the overfitting problem in decision trees for the training set. And the right kind of randomness makes them accurate classifiers and regressors. The output of the random forests for classification is the class chosen by the majority of the trees, and the regression task returns the mean prediction of the individual trees [41]. Excellent performance is shown by random forest when the number of variables is much larger than the number of observations; it allows hundreds of input variables without variable deletion. Furthermore, random forests have extremely high classification rates,

and they can make feature selection by determining the importance of variables [15]. All statistical analyses about random forest are using R Package *randomForest* [40]. The sample size of each bootstrap replicate is not the same as the original data, which is around 480 each time. There are 500 number of trees grown and 12 number of predictors sampled for splitting at each node.

Steps involved in random forest algorithm:

- Create a Bootstrapped Data Set.
- Individual decision trees are constructed for each sample.
- Each decision tree will generate an output.
- Go back to Step 1 and Repeat
- Bootstrapped the data and used the aggregate from all the trees to make a decision. This process is known as Bagging (Majority Voting or Averaging for Classification and regression, respectively).

Usually, there are two measures of variable importance for random forests: Mean Decrease Gini and Mean Decrease Accuracy. In this thesis, we use Mean Decrease Accuracy(%IncMSE), the most robust and informative measure obtained from RSS (for bagging regression trees). For calculating the %IncMSE for a given predictor, we record the total amount of RSS decrease due to splits over this given predictor and average overall bagging trees. The higher value of Mean Decrease Accuracy, the higher importance of the variable in the random forest model.

3.2.4 Predictive R^2 with CV

For some regularized methods (e.g. Lasso), we typically have numerous variables, and a model with all the variables will always have the largest R^2 . As such, we are using R^2 is not recommended. With a continuous dependent variable, we first divide the whole dataset into 10 folds and use one of the folds as a testing set and the others as a training set. Each fold will become a test set, making a total of ten predictions. Then we would use the Predicted Residual Sum of Squares (PRESS) statistic, which is given by

$$\text{PRESS} = \sum_{k=1}^K \sum_{i \in \text{fold}_k} (y_i - \hat{y}_{i,-k})^2 \quad (3.17)$$

where $\hat{y}_{i,-k}$ is the i^{th} predictive value for y_i without using test cases in fold_k . $k \in \{1, 2, 3, \dots, 10\}$ $i \in \text{fold}_k$.

The predictive R^2 is denoted by

$$R^2 = 1 - \frac{\text{PRESS}}{\text{SS}} \quad (3.18)$$

where $\text{SS} = \sum_i (y_i - \bar{y})^2$, \bar{y} is the mean of y_i in the whole dataset.

3.2.5 Cross Validation

Cross-validation is a statistical method used to evaluate the performance of a machine learning model. It is often used in applied machine learning to compare and select models for a given predictive modelling problem, as it helps us choose the model that performs best on unseen data, evaluate the quality of the model and avoid overfitting and underfitting. In this method, we split the data into train and test sets so that the sample between train and test set does not overlap[17]. As there is never enough data to train a model, removing a part of it for validation poses a problem of under-fitting. We risk losing essential patterns by reducing training data and increasing bias-induced error. K-Fold and Leave One Out Cross-Validation are popular and easy ways to understand. They generally result in a less biased model compared to other methods. Because it ensures that every observation from the original dataset has the chance of appearing in the training and test set, this is one of the best approaches if we have limited input data[16].

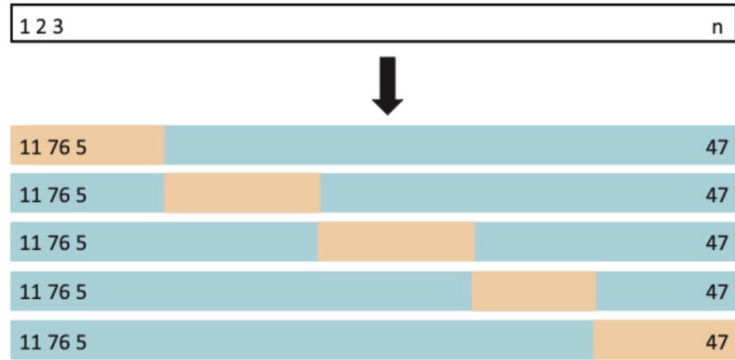


Figure 3.3: 5-fold Cross-Validation

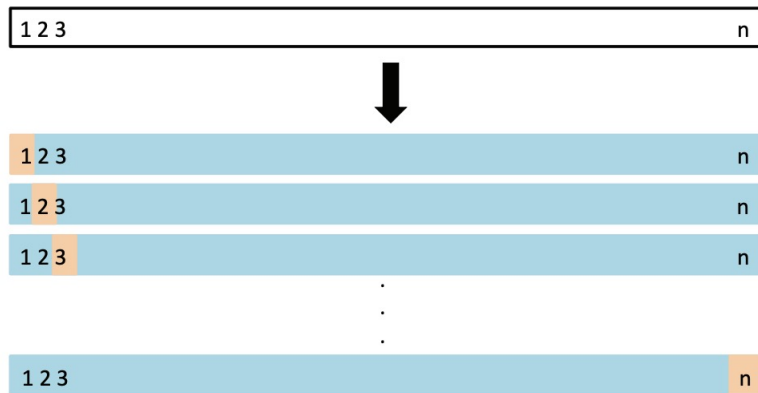


Figure 3.4: Leave One Out Cross-Validation

A set of n observations randomly divided the dataset into two segments. One segment is for training the model, and remains for testing the model. K-fold cross-validation means randomly partitioning all samples

into k equal-sized and non-overlapping parts. One of the k parts is regarded as the validation data for testing the model, and the remaining $k - 1$ parts are used as the training set. The k -fold cross-validation process repeats k times, with each k subset being used once as the validation set. Figure 3.1 is a schematic display of a 5 Fold cross-validation.

In LOOCV, we divide the data set into two parts. In one part, we have a single observation, which is our test data, and in the other part, we have all the other observations from the dataset forming our training data. If we have a data set with n observations, then training data contains $n - 1$ observation, and test data contains 1 observation. The first training set contains all but observation 1, the second training set contains all but observation 2, and so forth. This process is iterated for each data point as shown in Figure 3.2[16].

The test error is then estimated by averaging the n resulting MSEs. It is described as follows:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3.19)$$

where y_i is the observation and \hat{y}_i is the predictive value of this observation.

4 Result

4.1 FPCA

Elderly cognitive data contains one time-varying trait, the DRS. We obtained patterns of variation in the cognitive outcomes by using FPCA. Figure 4.1 shows the mean function, scree plot and the first three eigenfunctions. The y-axis in the mean function plot denotes the normalized DRS. The computed mean function for the data shows a dip in DRS near the beginning, then a clear upward trend at age 70, followed by an abrupt decline at age 75. In the scree plot, the y-axis represents FPCs, and the x-axis represents the fraction of variance explained. It can be seen that FPC1 explained a large fraction of the variation (78.0%). The second, third, and fourth FPCs explained 14.2, 6.7 and 1.1% of the variation, respectively. The first two FPC explained 92.2% of the variability, and 99.9% of the variability was described by the first 4 FPCs. The remaining are less important except for the first two FPCs. Eigenfunctions(EFs) are also known as weight functions in Chapter 3. Interpreting these EFs can be quite hard, as there might not be an obvious counterpart in the data. EFs of DRS seem to explain that EF1 is the reduction of the DRS from the mean function over age, EF2 is the reduction of DRS before and after the age of 80.

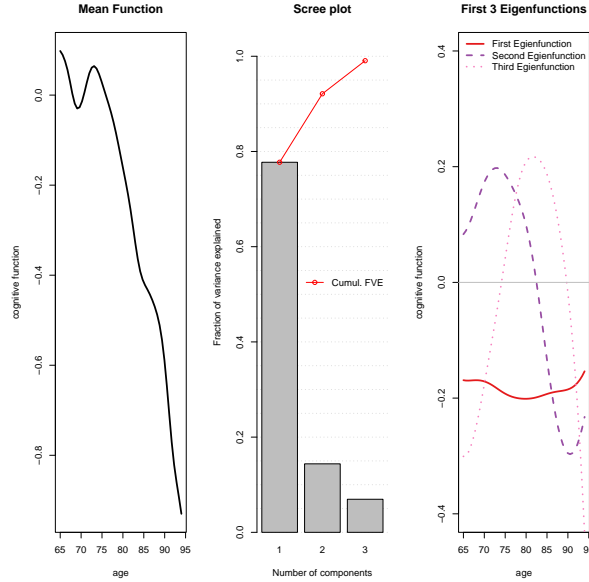


Figure 4.1: Summary Of FPCA Result

4.2 Comparison of FPC scores by countries

The sample data after filtering involved 9 European countries (Austria, Belgium, Denmark, France, Germany, Italy, Spain, Sweden, and Switzerland). We compare the first two FPCs and the median longitudinal trajectory for these nine countries in Figures 4.2 – 4.5. The curves in Figures 4.2 are the mean function plus the median FPC1 of each country. Since the EF1 is a decline in the cognitive mean function, We can regard these curves as the mean cognitive curve after the first dip. After adding the median FPC2 of each country, the curves shown in Figure 4.3 show another descent of the curves 4.2. It can be seen from both Figures that the median longitudinal trajectory curves (MLTC) of Italy and Spain are at the bottom. However, from Figure 4.4, we find that Italy and Spain have higher FPC1 than other counties, and Sweden and Switzerland have relatively high FPC1. This is because EF1 is negative. The first FPC for each country has a negative correlation with the actual DRS, which can be verified in Figures 4.2 and 4.3. We guess it is because these two countries' economies are not great compared to other European countries, resulting in an imperfect health care system and a lack of medical conditions.

Table 4.1 shows the Euro Health Consumer Index (EHCI) in 2015 (from the Euro health consumer index 2015 report), which was a comparison of European health care systems based on waiting times, results, and generosity and included 37 European countries. Based on patient rights and information score (PRIS), pharmaceuticals score and overall score ranking in the table, we find that Spain and Italy are in eighth and ninth positions for these indexes, respectively. In addition, these two countries do not have high scores on waiting times for treatment (scores) (WTTS), outcomes (scores) (OS), range and reach of services (scores) (RRSS) and prevention scores, which provide evidence for our results that the MLTC of these two countries is at the bottom among these nine European countries. Figure 4.1 shows that 80 is a turning point, and FPC2 will become negative after about 80 years old. According to the boxplot in Figure 4.5, only the medians of Spain and Italy are below zero. This result confirms the summary result mentioned in section 2.1, the severe problem of population ageing in Spain and Italy.

The top 5 MLTC are from Sweden, Switzerland, Denmark, Germany, and France. The comparative figure of the FPC2 (Figure 4.5) shows that the overall country difference in pattern change is not as substantial as the difference in the first FPC scores.

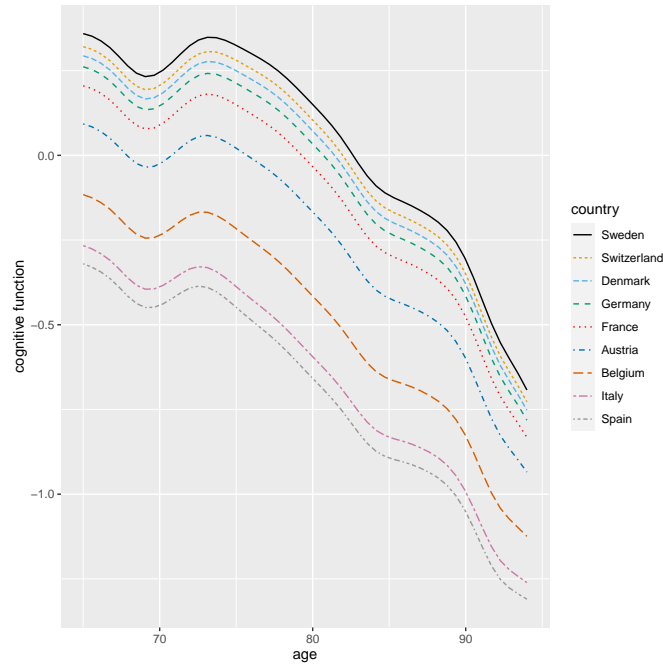


Figure 4.2: Overall Cognitive Function Trajectory Across Different European Countries Based On FPC1. (Mean function + Median FPC1 of each country)

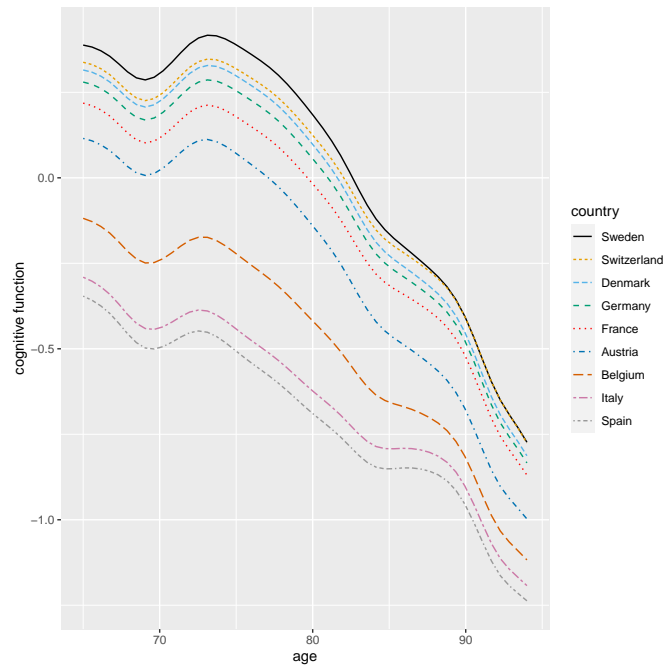


Figure 4.3: Overall Cognitive Function Trajectory Across Different European Countries Based On FPC1 and FPC2. (Mean function + Median FPC1 + Median FPC2 of each country)

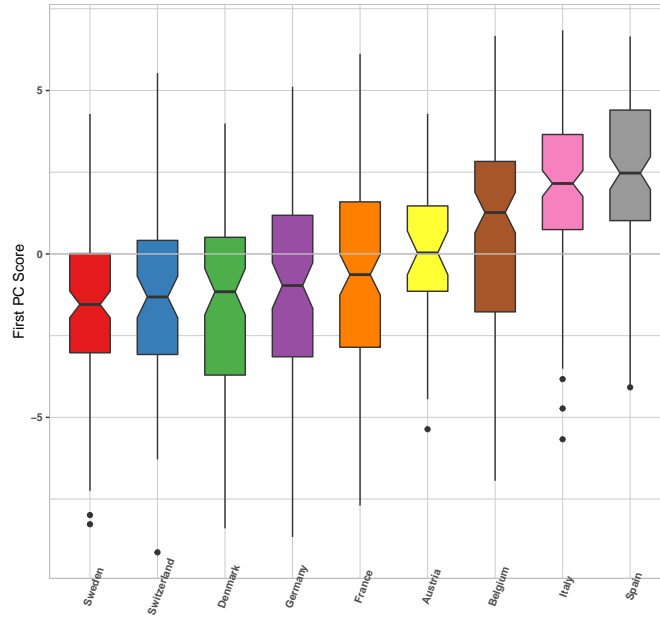


Figure 4.4: Bloxplot Of FPC1 Across Different European Countries

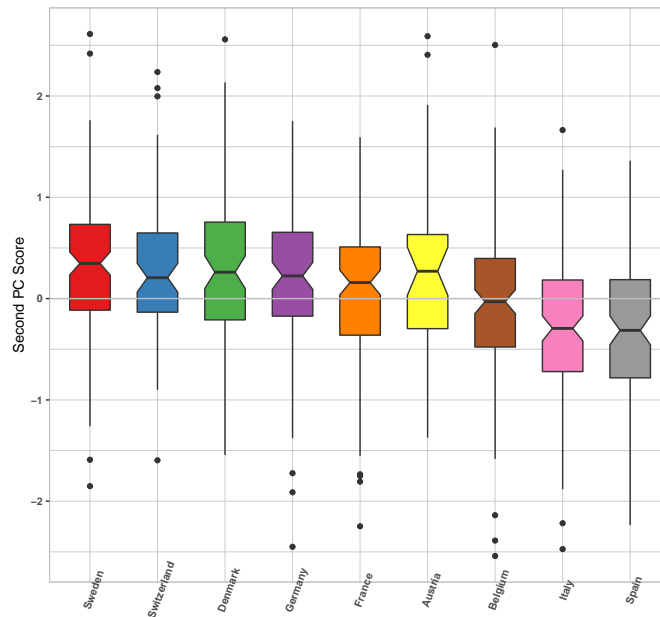


Figure 4.5: Bloxplot Of FPC2 Across Different European Countries

Table 4.1: Euro Health Consumer Index 2015

Country	Overall ranking	Total socre	PRIS	WTTS	OS	RRSS	Prevention score	PS
Switzerland	2	894	133	225	240	119	101	76
Belgium	5	836	117	225	198	131	89	76
Germany	7	828	125	188	229	94	107	86
Denmark	9	793	133	138	219	138	89	76
Sweden	10	786	125	100	229	144	107	81
France	11	775	113	188	208	106	89	71
Austria	12	774	121	188	188	119	83	76
Spain	19	695	104	113	198	113	101	67
Italy	22	667	96	138	188	88	101	57

4.3 Predictive R^2 for models

Table 4.2 shows the predictive R^2 of FPC1 and FPC2 for different models. We calculate predictive R^2 four times for each model (two with and two without covariate DRS). From the table, we conclude that the predictive R^2 for FPC1 with covariate DRS is around 0.4 to 0.5. If we remove the covariate DRS, the predictive R^2 drops from 0.4 to 0.2. Nevertheless, the Predictive R^2 of FPC2 for all models is very low, around 0.1, regardless of whether the covariate DRS is contained. The best result of predictive R^2 for predicting FPC1 is to use the Elastic-net model, which is 0.550 with covariate DRS and 0.450 without covariate DRS. Again, the best result of predictive R^2 for predicting FPC2 is to use the elastic-net model (0.132 with covariate DRS and 0.115 without covariate DRS). Since the predictive R^2 for FPC2 is low, we should not over-interpret the coefficients as we presented in the previous report. Many factors for explaining the reduction of cognitive functions after 80 years old haven't been included in this study, and this will be an exciting topic for future research.

Table 4.2: Predictive R^2 for machine learning models

Methods	For FPC1 with DRS	For FPC1 without DRS	For FPC2 with DRS	For FPC2 without DRS
Elastic-net	0.550	0.450	0.132	0.115
Lasso	0.540	0.400	0.130	0.080
Random Forest	0.479	0.370	0.122	0.100
Decision Tree	0.436	0.260	0.094	0.074

4.4 Indentification of Risk factors for cognitive decline

We set up a 10 fold cross-validation strategy for choosing the parameter α and λ in the elastic net model. And we use the argument “tuneLength” in R that tests different combinations of values for α and λ . There are ten α values from 0.10 to 1.00, and each tried ten times to select the optimal model. The smallest value

of the root mean square error (RMSE) is used to select the optimal model.

With FPC1 as dependent variable, the final tuning parameter used for the model are $\alpha = 0.2$ and $\lambda = 0.28$. Generally, a model with all the variables will have a larger R^2 . Hence for a regularized model, we prefer to use RMSE, mean absolute error(MAE) and predictive R^2 for goodness-of-fit measures. The $RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$ for the model include all covariates in Table 2.1 and 2.2 is 2.28 and $MAE = \frac{\sum_{i=1}^n |\hat{y}_i - y_i|}{n}$ is 1.87. If we include the covariate DRS the $RMSE$ and MAE decrease to 2.11 and 1.72, respectively. The predictive R^2 , including all covariates, is 0.42 and reduced to 0.55 if covariate DRS is added. With FPC2 as dependent variable, the final tuning parameter used for the model are $\alpha = 0.2$ and $\lambda = 0.28$. However, the predictive R^2 for FPC2 is low (0.132 with DRS and 0.115 without DRS). The details of the predictive R^2 are in Table 4.2.

In Figures 4.6 and 4.7, we report the risk factors of baseline covariates by using elastic-net. The positive, negative, and no effects with an order of importance are represented by red, blue, and grey, respectively. The top five positive effects in baseline wave for FPC1 are Gross Motor Skills level 2 out of 4, Country Spain, Country Italy, Fine Motor Skills level 2 out of 4 and Instrumental Activities of Daily level 1 out of 2. And top five negative effects for FPC1 are Immediate recall score, country Sweden, numeracy score of 5 out of 5, Country France and gender Female. Notice since EF1 is negative. All effects are opposite for the original DRS. Namely, A negative effect on FPC1 is a positive effect on DRS.

In general, the elderly living in countries with excellent health care systems, being in good health, having a long education experience, being accompanied by family members, exercising regularly, being in a good mood without depression, having an average cognitive ability higher than the average of all people. More specifically, We find that the Gross Motor Skills Index level 2 out of 4 is the factor that decreases the most of the mean cognitive function (Fig. 4.6). This Index describes the sum of walking 100 meters, walking across a room, climbing one flight of stairs and bathing or showering. The higher the index, the better the cognitive ability of the participant. Similarly, Fine Motor Skills Index is the sum of picking up a small coin, eating/cutting up food and dressing. Instrumental Activities are the sum of telephone calls, taking medications, and managing money. Fine Motor Skills Index level 2 and Instrumental Activities level 1 are relatively low activity indexes. These results illustrate the importance of activity for cognitive decline. In addition, as we mentioned in section 4.2, the country is a significant factor affecting cognitive function, and the same conclusion can be drawn from Figure 4.6 as in section 4.2. The elderly from Spain, Italy, and Belgium showed slower mean cognitive function decline, while the elderly from Sweden, France and Denmark showed a faster reduction. Furthermore, we found that the mean cognitive function decreased more in the elderly hospitalized and with a high depression index. Cognitive decline can be mitigated by long years of education, good reading, writing skills, and eyesight. Immediate recall score is similar to delay recall score (DRS), so it greatly inhibits mean cognitive function decline as we expected. Compared with the top effects in Figure 4.7, we find the significant effects in both figures are very similar. Some covariates are continually significant. (e.g., country, Gross Motor Skill, immediate recall score)

To validate and compare risk factors estimated by elastic-net shown in Figure 4.6-4.7, we report four graphs (see Fig. 4.8-4.9 in the main text and Figures in Appendix A) where we offer the tree-based method for predicting FPC1 and FPC2. We employ the tree’s CP value with the smallest cross-validation error. The decision tree algorithm can naturally select which features are most crucial [47]. Besides, we prune parts of the tree that do not provide the power to classify instances to reduce the probability of overfitting problems. The features that impact our outcome variable DRS demonstrated in these four figures are comparable to those in the elastic-net model. (e.g., immediate recall score, country, numeracy score, reading)

Nevertheless, decision trees still suffer from overfitting the training data. Thus, we also use a random forest model to identify significant variables (risk factors). As mentioned above, a random forest model is a way of producing multiple decision trees with different parts of the same training set, which is insensitive to the overfitting of an individual decision tree [4]. Mean Decrease Accuracy ranks the random forest algorithm’s essential features shown in Figure 4.13. Immediate recall score, country, numeracy score, reading and gender are still the top important variables.

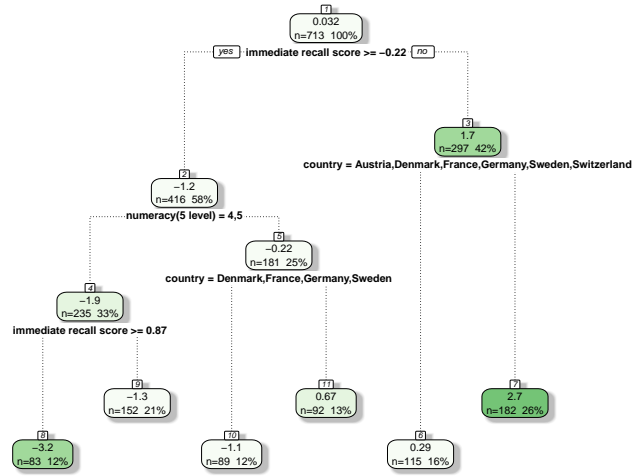


Figure 4.8: Pruned Tree based method for predicting FPC1 from all the variables

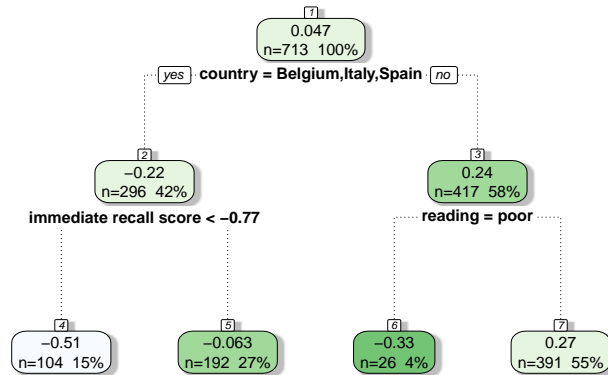


Figure 4.9: Pruned Tree based method for predicting FPC2 from all the variables

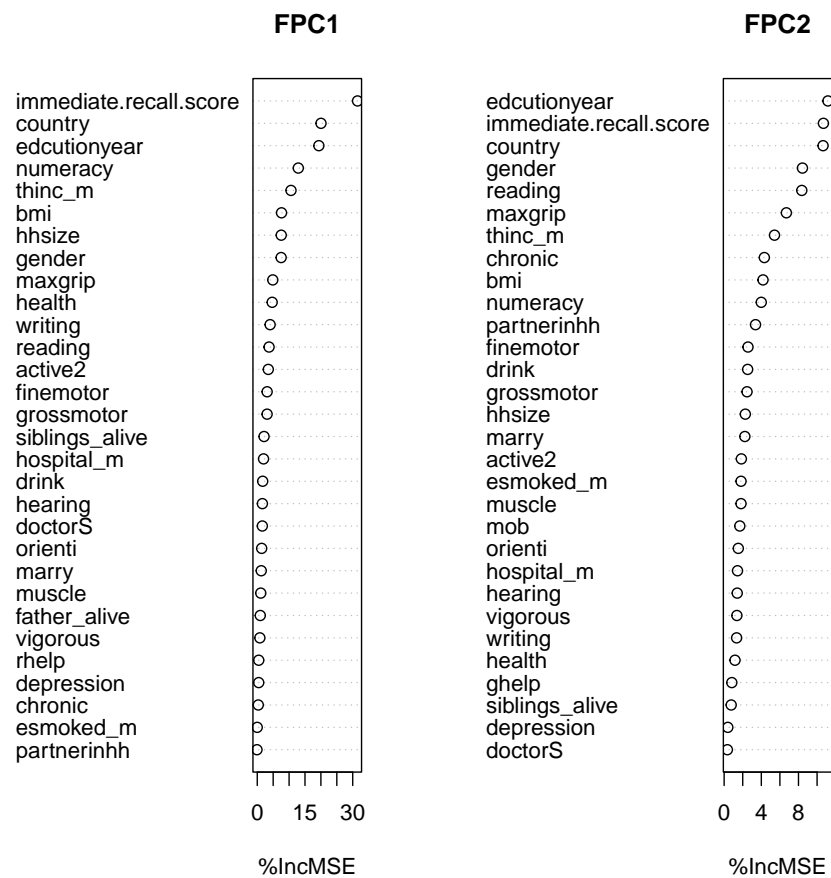


Figure 4.10: Top 30 ranked variables identified by Random Forest based on the Two FPC Scores

5 Conclusion

In this thesis, we have shown how the FPCA is helpful in the study of cognitive decline among age in the elderly, which can catch the major variance of the longitudinal cognitive measurements. Applying these functional techniques to European individuals aged 65 to 80 confirmed many previous conjectures and revealed many interesting findings. Our results emphasise the importance of education, Income, activities, drinking, and chronic diseases for cognitive decline. While some factors are unavoidable, others may not, such as lowering BMI through exercise or reducing alcohol consumption. Therefore, to prevent or delay cognitive decline in the elderly, interventions to prevent it may be an effective strategy.

Furthermore, according to our analysis, other European countries have similar cognitive abilities to people over 65, except for Italy and Spain. Nevertheless, they all go through two stages (early and late) of cognitive decline, and the decline will accelerate in the second stage (around 75). While the differences between Italy, Spain and other Europe countries are significant, further research is needed to discover the reason we find the link between country and cognitive decline. We believe this kind of analysis could be of considerable interest to cognitive patients since it allows them to develop a plan in accordance with early prevention. Based on the predictive R^2 we obtained, we believe that many other factors affect cognitive function in older adults, especially factors for explaining the reduction of cognitive functions after 80 years old, which is an exciting and worthwhile topic to explore in future research.

References

- [1] Liu B and Müller H G. Estimating derivatives for samples of sparsely observed functions, with application to online auction dynamics. *Journal of the American Statistical Association*, 2009.
- [2] Paul B Baltes, Ursula M Staudinger, and Ulman Lindenberger. Lifespan psychology: Theory and application to intellectual functioning. *Annual review of psychology*, 50:471–507, 1999.
- [3] Deborah E Barnes and Kristine Yaffe. The projected effect of risk factor reduction on alzheimer’s disease prevalence. *The Lancet Neurology*, 10(9):819–828, 2011.
- [4] Mariana Belgiu and Lucian Drăguț. Random forest in remote sensing: A review of applications and future directions. *ISPRS journal of photogrammetry and remote sensing*, 114:24–31, 2016.
- [5] Bennett and Derrick A. How can i deal with missing data in my study? *Australian and New Zealand journal of public health*, 25(5):464–469, 2001.
- [6] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 1. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [7] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 2. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [8] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 4. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [9] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 5. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [10] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 6. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [11] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 7. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [12] A. Börsch-Supan. Survey of health, ageing and retirement in europe (share) wave 8. release version: 8.0.0. *SHARE-ERIC. Data set.*, 2022.
- [13] Axel Börsch-Supan, Martina Brandt, Christian Hunkler, Thorsten Kneip, Julie Korbmacher, Frederic Malter, Barbara Schaan, Stephanie Stuck, and Sabrina Zuber. Data resource profile: the survey of health, ageing and retirement in europe (share). *International journal of epidemiology*, 42(4):992–1001, 2013.
- [14] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- [15] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [16] Michael W Browne. Cross-validation methods. *Journal of mathematical psychology*, 44(1):108–132, 2000.
- [17] Jason Brownlee. A gentle introduction to k-fold cross-validation. *Machine learning mastery*, 2019, 2018.
- [18] Edward Joseph Caruana, Marius Roman, Jules Hernández-Sánchez, and Piergiorgio Solli. Longitudinal studies. *Journal of thoracic disease*, 7(11):E537, 2015.

- [19] Kehui Chen, Xiaoke Zhang, Alexander Petersen, and Hans-Georg Müller. Quantifying infinite-dimensional data: Functional data analysis in action. *Statistics in Biosciences*, 9(2):582–604, 2017.
- [20] Castro P E, Lawton W H, and Sylvestre E A. Principal modes of variation for processes with continuous sample curves. *Technometrics*, 1986.
- [21] Yao F, Müller H G, and et al. Clifford A J. Shrinkage estimation for functional principal component scores with application to the population kinetics of plasma folate. *Biometrics*, 2003.
- [22] Jianqing Fan and Irene Gijbels. *Local polynomial modelling and its applications*. Routledge, 2018.
- [23] Frédéric Ferraty and Philippe Vieu. *Nonparametric functional data analysis: theory and practice*, volume 76. Springer, 2006.
- [24] Staniswalis J G and Lee J J. Nonparametric regression analysis of longitudinal data. *Journal of the American Statistical Association*, 1998.
- [25] David S Geldmacher and Peter J Whitehouse. Evaluation of dementia. *New England Journal of Medicine*, 335(5):330–336, 1996.
- [26] Peter Hall, Hans-Georg Müller, and Fang Yao. Modelling sparse generalized longitudinal observations with latent gaussian processes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(4):703–723, 2008.
- [27] Kyunghye Han, Pantelis Z Hadjipantelis, Jane-Ling Wang, Michael S Kramer, Seungmi Yang, Richard M Martin, and Hans-Georg Müller. Functional principal component analysis for identifying multivariate patterns and archetypes of growth, and their association with long-term cognitive development. *PloS one*, 13(11):e0207073, 2018.
- [28] Lajos Horváth and Piotr Kokoszka. *Inference for functional data with applications*, volume 200. Springer Science & Business Media, 2012.
- [29] Fan J and Gijbels I. Local polynomial modelling and its applications. *Routledge*, 2018.
- [30] Reka Karuppusami, Belavendra Antonisamy, and Prasanna S Premkumar. Functional principal component analysis for identifying the child growth pattern using longitudinal birth cohort data. *BMC Medical Research Methodology*, 22(1):1–10, 2022.
- [31] Kim Kiely. Cognitive function. *Encyclopedia of Quality of Life and Well-Being Research*. Springer, 2014.
- [32] Max Kuhn, Jed Wing, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, Brenton Kenkel, R Core Team, et al. Package ‘caret’. *The R Journal*, 22:7, 2020.
- [33] Paul Madley-Dowd, Rachael Hughes, Kate Tilling, and Jon Heron. The proportion of missing data should not be used to guide decisions on multiple imputation. *Journal of clinical epidemiology*, 110:63–73, 2019.
- [34] Petra Maresova, Blanka Klimova, Michal Novotny, and Kamil Kuca. Alzheimer’s and parkinson’s diseases: Expected economic impact on europe—a call for a uniform european strategy. *Journal of Alzheimer’s Disease*, 54(3):1123–1133, 2016.
- [35] Stephen Milborrow and Maintainer Stephen Milborrow. Package ‘rpart. plot’. *Plot’rpart’Models: An Enhanced Version of rpart*, 2019.
- [36] Hans-Georg Müller. Peter hall, functional data analysis and random objects. *The Annals of Statistics*, 44(5):1867–1887, 2016.
- [37] DL Murman. The impact of age on cognition. In *Seminars in Hearing*, volume 36, pages 111–121, 2015.
- [38] HI Niu, F Álvarez-Álvarez, and I Guillén-Grima. Aguinaga-ontoso. prevalence and incidence of alzheimer’s disease in europe: A meta-analysis. *Neurología*, 32(8):523–532, 2017.

- [39] Sam Norton, Fiona E Matthews, Deborah E Barnes, Kristine Yaffe, and Carol Brayne. Potential for primary prevention of alzheimer’s disease: an analysis of population-based data. *The Lancet Neurology*, 13(8):788–794, 2014.
- [40] Suggests RColorBrewer and Maintainer Andy Liaw. Package ‘randomforest’. *University of California, Berkeley: Berkeley, CA, USA*, 2018.
- [41] Victor Francisco Rodriguez-Galiano, Bardan Ghimire, John Rogan, Mario Chica-Olmo, and Juan Pedro Rigol-Sanchez. An assessment of the effectiveness of a random forest classifier for land-cover classification. *ISPRS journal of photogrammetry and remote sensing*, 67:93–104, 2012.
- [42] Donald B Rubin. Multiple imputation for nonresponse in surveys. hoboken, 1987.
- [43] Joseph L Schafer. Multiple imputation: a primer. *Statistical methods in medical research*, 8(1):3–15, 1999.
- [44] Haolun Shi, Jianghu Dong, Liangliang Wang, and Jiguo Cao. Functional principal component analysis for longitudinal data with informative dropout. *Statistics in Medicine*, 40(3):712–724, 2021.
- [45] Joan G Staniswalis and J Jack Lee. Nonparametric regression analysis of longitudinal data. *Journal of the American Statistical Association*, 93(444):1403–1418, 1998.
- [46] Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics*, 8(1):1–21, 2007.
- [47] Steinbach M TAN PN and Vipin Kumar. Introduction to data mining, 2014.
- [48] Terry Therneau, Beth Atkinson, Brian Ripley, and Maintainer Brian Ripley. Package ‘rpart’. *Available online: cran. ma. ic. ac. uk/web/packages/rpart/rpart. pdf (accessed on 20 April 2016)*, 2015.
- [49] Stef Van Buuren, Jaap PL Brand, Catharina GM Groothuis-Oudshoorn, and Donald B Rubin. Fully conditional specification in multivariate imputation. *Journal of statistical computation and simulation*, 76(12):1049–1064, 2006.
- [50] Anders Wimo, Linus Jönsson, John Bond, Martin Prince, Bengt Winblad, and Alzheimer Disease International. The worldwide economic impact of dementia 2010. *Alzheimer’s & dementia*, 9(1):1–11, 2013.
- [51] Bengt Winblad, Philippe Amouyel, Sandrine Andrieu, Clive Ballard, Carol Brayne, Henry Brodaty, Angel Cedazo-Minguez, Bruno Dubois, David Edvardsson, Howard Feldman, et al. Defeating alzheimer’s disease and other dementias: a priority for european science and society. *The Lancet Neurology*, 15(5):455–532, 2016.
- [52] Brita Askeland Winje, Jo Røislien, Eli Saastad, Jorid Eide, Christopher Finne Riley, Babill Stray-Pedersen, and J Frederik Frøen. Wavelet principal component analysis of fetal movement counting data preceding hospital examinations due to decreased fetal movement: a prospective cohort study. *BMC pregnancy and childbirth*, 13(1):1–11, 2013.
- [53] Hans-Ulrich Wittchen, Frank Jacobi, Jürgen Rehm, Anders Gustavsson, Mikael Svensson, Bengt Jönsson, Jes Olesen, Christer Allgulander, Jordi Alonso, Carlo Faravelli, et al. The size and burden of mental disorders and other disorders of the brain in europe 2010. *European neuropsychopharmacology*, 21(9):655–679, 2011.
- [54] Fang Yao, Hans-Georg Müller, and Jane-Ling Wang. Functional data analysis for sparse longitudinal data. *Journal of the American statistical association*, 100(470):577–590, 2005.
- [55] John Zeisel, Kirsty Bennett, and Richard Fleming. World alzheimer report 2020: Design, dignity, dementia: Dementia-related design and the built environment. 2020.
- [56] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2):301–320, 2005.

Appendix A

Plots for predicting the first two FPC scores from the categorical variables

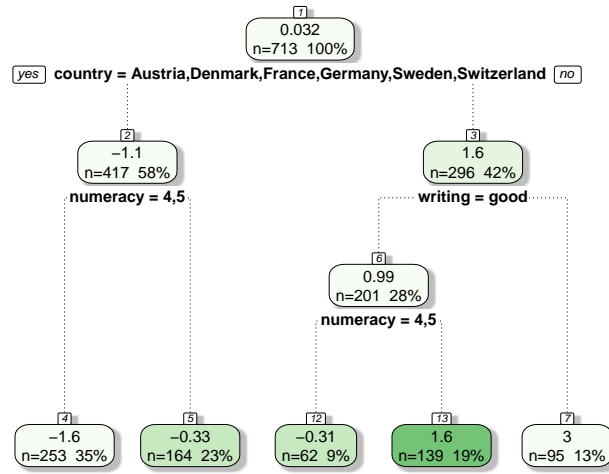


Figure A.1: Pruned Tree based method for predicting FPC1 from the categorical variables

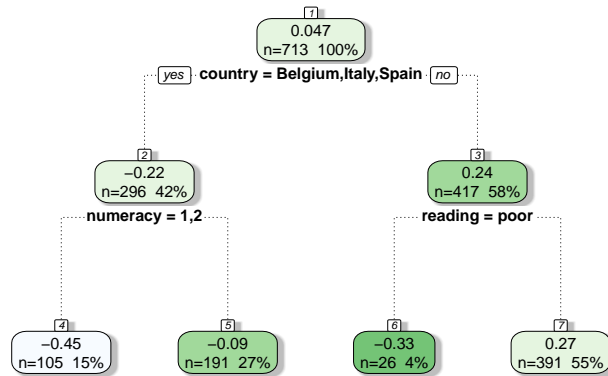


Figure A.2: Pruned Tree based method for predicting FPC2 from the categorical variables

Appendix B

R Code

B.1 Data Management for Wave8

B.2 Load data and data clearing

```
1 #Note this is only the data management for wave 8, the data management
  for wave1-7 data is similar to this procedure.
2
3 #clean function
4 cleanup <- function(myvar = Alcohol$ALQ110, mymin = 1, mymax = 2,
  exclude = c(999)){
5   ID1 <- (myvar < mymin)           # ids with less than min values
6   ID2 <- (myvar > mymax)           # ids with greater than max values
7   if (length(exclude) > 0)        # ids with strange values
8   { ID3 <- myvar%in%exclude }
9
10  # set all the above three values to be NA.
11  RmID <- ((ID1 + ID2 + ID3) > 0)
12  myvar[RmID] = NA
13  return(myvar)
14 }
15
16 data=read.csv("wave1_7.csv")
17 library(haven)
18 cr <- read_dta('sharew8_rel8-0-0_cv_r.dta')
19 cf <- read_dta('sharew8_rel8-0-0_cf.dta')
20 gv_health<-read_dta('sharew8_rel8-0-0_gv_health.dta')
21 names(cr)[names(cr)=="hhid8"]="hhid"
22 names(cr)[names(cr)=="mergeidp8"]="mergeidp"
23 names(cr)[names(cr)=="coupleid8"]="coupleid"
24 names(gv_health)[names(gv_health)=="hhid8"]="hhid"
25 names(gv_health)[names(gv_health)=="mergeidp8"]="mergeidp"
26 names(gv_health)[names(gv_health)=="coupleid8"]="coupleid"
27 names(gv_health)[names(gv_health)=="cf008tot"]="recall_1"
28 names(gv_health)[names(gv_health)=="cf016tot"]="recall_2"
29 cr$wave <- 8
30 gv_health$wave <- 8
31 inter<-c(intersect(colnames(cr), colnames(data)))
32 inter2<-c(intersect(colnames(gv_health), colnames(data)))
33
34 wave8cr<-cr[,c(inter)]
35 wave8gv_health<-gv_health[,c(inter2)]
36 wave81 <- merge(wave8cr, wave8gv_health, by = intersect(names(wave8cr),
  names(wave8gv_health)), all = TRUE)
37
38 ac=read_dta('sharew8_rel8-0-0_ac.dta')
39 names(ac)[names(ac)=="hhid8"]="hhid"
```

```

40 names(ac)[names(ac)=="mergeidp8"]=="mergeidp"
41 names(ac)[names(ac)=="coupleid8"]=="coupleid"
42 ac$wave=8
43 ac$ac002d1=NA
44 ac$ac002d2=NA
45 ac$ac002d3=NA
46 ac$ac002d4=NA
47 ac$ac002d5=NA
48 ac$ac002d6=NA
49 ac$ac002d7=NA
50 ac$ac002dno=NA
51 selectac=c("mergeid","hhid","coupleid","wave","ac002d1","ac002d2","
      ac002d3","ac002d4","ac002d5",
52      "ac002d6","ac002d7","ac002dno")
53 wave8ac<-ac[,selectac]
54 wave81new <- merge(wave81,wave8ac,by = intersect(names(wave81), names(
      wave8ac)), all = TRUE)
55 #write.csv(wave81,file = "Desktop/2022\ share/wave81.csv",row.names = F)
56 #####
57 #####part 2
58 dn<-read_dta('sharew8_rel8-0-0-dn.dta')
59 names(dn)[names(dn)=="hhid8"]=="hhid"
60 names(dn)[names(dn)=="mergeidp8"]=="mergeidp"
61 names(dn)[names(dn)=="coupleid8"]=="coupleid"
62 names(dn)[names(dn)=="dn014_"]=="mar_stat"
63 names(dn)[names(dn)=="dn026_1"]=="mother_alive"
64 names(dn)[names(dn)=="dn026_2"]=="father_alive"
65
66 dn <- dn %>% rowwise() %>%
67   mutate(siblings_alive = sum(dn036_,dn037_))
68 dn$wave<-8
69 selectdn=c("mergeid","hhid","coupleid","dn002_","dn003_","dn004_","
70      "dn007_","mar_stat","mother_alive","father_alive","dn037_","
      siblings_alive","wave")
71 wave8dn<-dn[,selectdn]
72
73 ch<-read_dta('sharew8_rel8-0-0-ch.dta')
74 names(ch)[names(ch)=="hhid8"]=="hhid"
75 names(ch)[names(ch)=="mergeidp8"]=="mergeidp"
76 names(ch)[names(ch)=="coupleid8"]=="coupleid"
77
78 ch$ch007_1[is.na(ch$ch007_1)]=0
79 ch$ch007_2[is.na(ch$ch007_2)]=0
80 ch$ch007_3[is.na(ch$ch007_3)]=0
81 ch$ch007_4[is.na(ch$ch007_4)]=0
82 ch$ch007_5[is.na(ch$ch007_5)]=0
83 ch$ch007_6[is.na(ch$ch007_6)]=0
84 ch$ch007_7[is.na(ch$ch007_7)]=0
85 ch$ch007_8[is.na(ch$ch007_8)]=0
86 ch$ch007_9[is.na(ch$ch007_9)]=0
87 ch$ch007_10[is.na(ch$ch007_10)]=0
88 ch$ch007_11[is.na(ch$ch007_11)]=0
89 ch$ch007_12[is.na(ch$ch007_12)]=0
90 ch$ch007_13[is.na(ch$ch007_13)]=0

```

```

91 ch$ch007_14[is.na(ch$ch007_14)]=0
92 ch$ch007_15[is.na(ch$ch007_15)]=0
93 ch$ch007_16[is.na(ch$ch007_16)]=0
94
95 ch$ch007_hh <- ifelse(ch$ch007_1 ==1| ch$ch007_2==1
96                       | ch$ch007_3==1| ch$ch007_4==1
97                       | ch$ch007_5==1| ch$ch007_6==1
98                       | ch$ch007_7==1| ch$ch007_8==1
99                       | ch$ch007_9==1| ch$ch007_10==1
100                      |ch$ch007_11==1|ch$ch007_12==1
101                      |ch$ch007_13==1|ch$ch007_14==1
102                      |ch$ch007_15==1|ch$ch007_16==1|
103                      ch$ch007_1 ==2| ch$ch007_2==2
104                      | ch$ch007_3==2| ch$ch007_4==2
105                      | ch$ch007_5==2| ch$ch007_6==2
106                      | ch$ch007_7==2| ch$ch007_8==2
107                      | ch$ch007_9==2| ch$ch007_10==2
108                      |ch$ch007_11==2|ch$ch007_12==2
109                      |ch$ch007_13==2|ch$ch007_14==2
110                      |ch$ch007_15==2|ch$ch007_16==2,1,
111                      ifelse(ch$ch007_1==0&ch$ch007_2==0
112                             &ch$ch007_3==0&ch$ch007_4==0
113                             &ch$ch007_5==0&ch$ch007_6==0
114                             &ch$ch007_7==0&ch$ch007_8==0&
115                             ch$ch007_9==0&ch$ch007_10==0&
116                             ch$ch007_11==0&ch$ch007_12==0&
117                             ch$ch007_13==0&ch$ch007_14==0&
118                             ch$ch007_15==0&ch$ch007_16==0,NA,5) )
119 table(ch$ch007_hh,useNA = "always")
120
121 ch$ch007_km <- ifelse(ch$ch007_1 ==3| ch$ch007_2==3
122                       | ch$ch007_3==3| ch$ch007_4==3
123                       | ch$ch007_5==3| ch$ch007_6==3
124                       | ch$ch007_7==3| ch$ch007_8==3
125                       | ch$ch007_9==3| ch$ch007_10==3
126                      |ch$ch007_11==3|ch$ch007_12==3
127                      |ch$ch007_13==3|ch$ch007_14==3
128                      |ch$ch007_15==3|ch$ch007_16==3,1,
129                      ifelse(ch$ch007_1==0&ch$ch007_2==0
130                             &ch$ch007_3==0&ch$ch007_4==0
131                             &ch$ch007_5==0&ch$ch007_6==0
132                             &ch$ch007_7==0&ch$ch007_8==0&
133                             ch$ch007_9==0&ch$ch007_10==0&
134                             ch$ch007_11==0&ch$ch007_12==0&
135                             ch$ch007_13==0&ch$ch007_14==0&
136                             ch$ch007_15==0&ch$ch007_16==0,NA,5) )
137 table(ch$ch007_km)
138
139 inter<-c(intersect(colnames(dn), colnames(data)))
140 inter
141 inter2<-c(intersect(colnames(ch), colnames(data)))
142 inter2
143
144 wave8dn<-dn[,c(inter)]

```



```

145 wave8ch<-ch[,c(inter2)]
146 wave82 <- merge(wave8dn, wave8ch, by = intersect(names(wave8dn), names(
      wave8ch)), all = TRUE)
147 #write.csv(wave82, file = "Desktop/2022\ share/wave82.csv", row.names = F)
148
149 #####
150 #####part 3
151 sp<-read_dta('sharew8_rel8-0-0_sp.dta')
152 names(sp)[names(sp)=="hhid8"]="hhid"
153 names(sp)[names(sp)=="mergeidp8"]="mergeidp"
154 names(sp)[names(sp)=="coupleid8"]="coupleid"
155 sp$wave<-8
156 selectsp<-c("mergeid", "hhid", "coupleid", "wave", "sp002_", "sp003_1", "sp003
      _2", "sp003_3", "sp008_",
157             "sp009_1", "sp009_2", "sp009_3")
158 wave8sp<-sp[,selectsp]
159
160
161 ph <-read_dta('sharew8_rel8-0-0_ph.dta')
162 ph$wave=8
163 names(ph)[names(ph)=="hhid8"]="hhid"
164 names(ph)[names(ph)=="mergeidp8"]="mergeidp"
165 names(ph)[names(ph)=="coupleid8"]="coupleid"
166 selectph=c("mergeid", "hhid", "coupleid", "wave", 'ph003_', paste0("ph006d"
      ,1:6), paste0("ph006d",10:14))
167 wave8ph<-ph[,selectph]
168
169
170 mh <-read_dta('sharew8_rel8-0-0_mh.dta')
171 mh$wave=8
172 names(mh)[names(mh)=="hhid8"]="hhid"
173 names(mh)[names(mh)=="mergeidp8"]="mergeidp"
174 names(mh)[names(mh)=="coupleid8"]="coupleid"
175 selectmh=c("mergeid", "hhid", "coupleid", "wave", "mh002_",
176            "mh003_", "mh004_", "mh005_", "mh007_", "mh008_", 'mh010_', "mh011_
      ",
177            "mh013_", "mh014_", "mh015_", "mh016_", "mh017_")
178
179 wave8mh<-mh[,selectmh]
180 wave8mh$mh002_<-ifelse(wave8mh$mh002_==1,1,ifelse(wave8mh$mh002_==5,0,NA
      ))
181 wave8mh$mh003_<-ifelse(wave8mh$mh003_==1,0,ifelse(wave8mh$mh003_==2,1,NA
      ))
182 wave8mh$mh004_<-ifelse(wave8mh$mh004_==1,1,ifelse(wave8mh$mh004_==2,0,NA
      ))
183 wave8mh$mh005_<-ifelse(wave8mh$mh005_==2,0,ifelse(wave8mh$mh005_==1,1,
      ifelse(wave8mh$mh005_==3,1,NA)))
184
185 wave8mh$mh007_<-ifelse(wave8mh$mh007_==1,1,ifelse(wave8mh$mh007_==2,0,NA
      ))
186 wave8mh$mh008_<-ifelse(wave8mh$mh008_==2,0,ifelse(wave8mh$mh008_==1,1,
      ifelse(wave8mh$mh008_==3,1,NA)))
187 wave8mh$mh010_<-ifelse(wave8mh$mh010_==1,1,ifelse(wave8mh$mh010_==2,0,NA
      ))

```

```

188 wave8mh$mh011_<-ifelse (wave8mh$mh011_==1,1,ifelse (wave8mh$mh011_==2,0,
    ifelse (wave8mh$mh011_==3,0,NA)))
189 wave8mh$mh013_<-ifelse (wave8mh$mh013_==1,1,ifelse (wave8mh$mh013_==5,0,NA
    ))
190 wave8mh$mh014_<-ifelse (wave8mh$mh014_==1,1,ifelse (wave8mh$mh014_==2,0,NA
    ))
191 wave8mh$mh015_<-ifelse (wave8mh$mh015_==1,1,ifelse (wave8mh$mh015_==2,0,NA
    ))
192 wave8mh$euro10=ifelse (wave8mh$mh014_==1|wave8mh$mh015_==1,1,0)
193 wave8mh$mh016_<-ifelse (wave8mh$mh016_==1,1,ifelse (wave8mh$mh016_==2,0,NA
    ))
194 wave8mh$mh017_<-ifelse (wave8mh$mh017_==1,1,ifelse (wave8mh$mh017_==5,0,NA
    ))
195
196 names (wave8mh) [names (wave8mh)=="mh002_"]="euro1"
197 names (wave8mh) [names (wave8mh)=="mh003_"]="euro2"
198 names (wave8mh) [names (wave8mh)=="mh004_"]="euro3"
199 names (wave8mh) [names (wave8mh)=="mh005_"]="euro4"
200 names (wave8mh) [names (wave8mh)=="mh007_"]="euro5"
201 names (wave8mh) [names (wave8mh)=="mh008_"]="euro6"
202 names (wave8mh) [names (wave8mh)=="mh010_"]="euro7"
203 names (wave8mh) [names (wave8mh)=="mh011_"]="euro8"
204 names (wave8mh) [names (wave8mh)=="mh013_"]="euro9"
205 names (wave8mh) [names (wave8mh)=="mh016_"]="euro11"
206 names (wave8mh) [names (wave8mh)=="mh017_"]="euro12"
207
208 wave8mh=subset (wave8mh,select = -c(mh014_,mh015_))
209 library (dplyr)
210 wave8mh <- wave8mh %>% rowwise() %>%
211   mutate(eurod = sum(c_>across(euro1:euro10)))
212
213 hc<-read_dta( 'sharew8_rel8-0-0-hc.dta' )
214 hc$wave=8
215 names (hc) [names (hc)=="hhid8"]="hhid"
216 names (hc) [names (hc)=="mergeidp8"]="mergeidp"
217 names (hc) [names (hc)=="coupleid8"]="coupleid"
218 names (hc) [names (hc)=="hc602_"]="hc002_"
219 selecthc<-c("mergeid","hhid","coupleid","wave","hc002_","hc012_","hc029_"
    ")
220 wave8hc<-hc[,selecthc]
221
222 ph<-read_dta( 'sharew8_rel8-0-0-ph.dta' )
223 ph$wave=8
224 names (ph) [names (ph)=="hhid8"]="hhid"
225 names (ph) [names (ph)=="mergeidp8"]="mergeidp"
226 names (ph) [names (ph)=="coupleid8"]="coupleid"
227 selectph=c("mergeid","hhid","coupleid","wave","ph049d1","ph049d2","
    ph049d3","ph049d4","ph049d5",
228   "ph049d8","ph049d9","ph049d10","ph049d11","ph049d13","ph048d1
    ","ph048d2","ph048d3","ph048d4","ph048d5",
229   "ph048d6","ph048d8","ph048d10")
230 wave8ph<-ph[,selectph]
231 wave8ph$ph049d1=cleanup (wave8ph$ph049d1,0,1)
232 wave8ph$ph049d3=cleanup (wave8ph$ph049d3,0,1)

```

```

233 wave8ph$ph049d4=cleanup (wave8ph$ph049d4,0,1)
234 wave8ph$ph049d2=cleanup (wave8ph$ph049d2,0,1)
235 wave8ph$ph049d5=cleanup (wave8ph$ph049d5,0,1)
236 wave8ph$ph049d10=cleanup (wave8ph$ph049d10,0,1)
237 wave8ph$ph049d11=cleanup (wave8ph$ph049d11,0,1)
238 wave8ph$ph049d13=cleanup (wave8ph$ph049d13,0,1)
239 wave8ph$ph049d9=cleanup (wave8ph$ph049d9,0,1)
240 wave8ph$ph049d8=cleanup (wave8ph$ph049d8,0,1)
241 wave8ph$ph048d1=cleanup (wave8ph$ph048d1,0,1)
242 wave8ph$ph048d4=cleanup (wave8ph$ph048d4,0,1)
243 wave8ph$ph048d5=cleanup (wave8ph$ph048d5,0,1)
244 wave8ph$ph048d2=cleanup (wave8ph$ph048d2,0,1)
245 wave8ph$ph048d3=cleanup (wave8ph$ph048d3,0,1)
246 wave8ph$ph048d6=cleanup (wave8ph$ph048d6,0,1)
247 wave8ph$ph048d8=cleanup (wave8ph$ph048d8,0,1)
248 wave8ph$ph048d10=cleanup (wave8ph$ph048d10,0,1)
249
250 library(dplyr)
251 wave8ph <- wave8ph %>% rowwise() %>%
252   mutate(adlwa = sum(c_across(ph049d1:ph049d4)))
253 table(wave8ph$adlwa,useNA = "always")
254
255 wave8ph <- wave8ph %>% rowwise() %>%
256   mutate(adla = sum(c_across(ph049d1:ph049d5)))
257
258
259 wave8ph <- wave8ph %>% rowwise() %>%
260   mutate(iadla = sum(c_across(ph049d10:ph049d13)))
261
262
263 wave8ph <- wave8ph %>% rowwise() %>%
264   mutate(iadlza = sum(c_across(ph049d8:ph049d10)))
265
266 wave8ph <- wave8ph %>% rowwise() %>%
267   mutate(mobilityind = sum(ph048d1,ph049d2,ph048d4,ph048d5))
268 table(wave8ph$mobilityind,useNA = "always")
269
270
271 wave8ph <- wave8ph %>% rowwise() %>%
272   mutate(lgmuscle = sum(ph048d2,ph048d3,ph048d6,ph048d8))
273
274
275 wave8ph <- wave8ph %>% rowwise() %>%
276   mutate(grossmotor = sum(ph048d1,ph049d2,ph048d5,ph049d3))
277
278
279 wave8ph <- wave8ph %>% rowwise() %>%
280   mutate(finemotor = sum(ph048d10,ph049d4,ph049d1))
281
282 select=c("mergeid","hhid","coupleid","wave","adlwa","adla","iadla",
283         "iadlza","mobilityind","lgmuscle","grossmotor",'finemotor')
284
285 wave8ph=wave8ph[,select]
286

```

```

287 M1<-merge(wave8hc, wave8sp, by = intersect(names(wave8hc), names(wave8sp))
      , all = TRUE)
288 M2=merge(M1, wave8ph, by = intersect(names(M1), names(wave8ph)), all = TRUE
      )
289 wave83=merge(M2, wave8mh, by = intersect(names(M2), names(wave8mh)), all =
      TRUE)
290
291 wave83$hc002_ =cleanup(wave83$hc002_, 0, 365)
292 wave83$hc012_ =cleanup(wave83$hc012_, 0, 5)
293 wave83$hc029_ =cleanup(wave83$hc029_, 1, 5)
294 wave83$sp002_ =cleanup(wave83$sp002_, 1, 5)
295 wave83$sp003_1=cleanup(wave83$sp003_1, 1, 96)
296 wave83$sp008_ =cleanup(wave83$sp008_, 1, 5)
297 wave83$sp009_1=cleanup(wave83$sp009_1, 1, 96)
298 wave83$sp009_2=cleanup(wave83$sp009_2, 1, 96)
299
300 inter<-c(intersect(colnames(wave83), colnames(data)))
301 inter
302 #write.csv(wave83, file = "wave83.csv", row.names = F)
303
304 #####
305 #####Part 4
306 br<-read_dta('sharew8_rel8-0-0_br.dta')
307 names(br)[names(br)=="hhid8"]="hhid"
308 names(br)[names(br)=="mergeidp8"]="mergeidp"
309 names(br)[names(br)=="coupleid8"]="coupleid"
310 br$wave <- 8
311 br$br010_mod<-NA
312 names(br)[names(br)=="br002_"]="smoking"
313 names(br)[names(br)=="br001_"]="ever_smoked"
314 selectbr<-c("mergeid", "hhid", "coupleid", "wave",
315             "smoking", "ever_smoked", "br010_mod", "br015_")
316 wave8br<-br[, selectbr]
317 wave83new=merge(wave83, wave8br, by = intersect(names(wave83), names(
      wave8br)), all = TRUE)
318 #write.csv(wave83new, file = "wave83.csv", row.names = F)
319
320 #####
321 #####Part 5
322 ep<-read_dta('sharew8_rel8-0-0_ep.dta')
323 names(ep)[names(ep)=="hhid8"]="hhid"
324 names(ep)[names(ep)=="mergeidp8"]="mergeidp"
325 names(ep)[names(ep)=="coupleid8"]="coupleid"
326 ep$wave <- 8
327 ep$ep011_ =NA
328 selectep<-c("mergeid", "hhid", "coupleid",
329             "wave", "ep005_", "ep009_", "ep011_", "ep013_", "ep026_", "ep036_"
330             )
331 wave8ep<-ep[, selectep]
332 wave83new2=merge(wave83new, wave8ep, by = intersect(names(wave83new),
      names(wave8ep)), all = TRUE)
333 #library(sjlabelled)
334 #a=as.list(get_label(ep))
335 #write.csv(wave83new2, file = "wave83.csv", row.names = F)

```

```

335
336
337 ph<-read_dta('sharew8_rel8-0-0_ph.dta')
338 ph$wave=8
339 names(ph)[names(ph)=="hhid8"]="hhid"
340 names(ph)[names(ph)=="mergeidp8"]="mergeidp"
341 names(ph)[names(ph)=="coupleid8"]="coupleid"
342 inter=c(intersect(colnames(ph), colnames(data)))
343 inter
344 wave8phnew=ph[,inter]
345 wave8phnew <- wave8phnew %>% rowwise() %>%
346   mutate(chronic_mod = sum(c_across(ph006d1:ph006d14)))
347
348 wave83new3=merge(wave83new2, wave8phnew, by = intersect(names(wave83new2),
349   names(wave8phnew)), all = TRUE)
349 #write.csv(wave83new3, file = "wave83.csv", row.names = F)
350
351
352 #####
353 #####Part 6
354
355 co<-read_dta('sharew8_rel8-0-0_co.dta')
356 co$wave=8
357 names(co)[names(co)=="hhid8"]="hhid"
358 names(co)[names(co)=="mergeidp8"]="mergeidp"
359 names(co)[names(co)=="coupleid8"]="coupleid"
360 selectco<-c("mergeid", "hhid", "coupleid", "wave", "co007_")
361 wave8co<-co[,selectco]
362
363 gv_isced<-read_dta('sharew8_rel8-0-0_gv_isced.dta')
364 names(gv_isced)[names(gv_isced)=="hhid8"]="hhid"
365 names(gv_isced)[names(gv_isced)=="mergeidp8"]="mergeidp"
366 names(gv_isced)[names(gv_isced)=="coupleid8"]="coupleid"
367 gv_isced$wave<-8
368 selectr<-c("mergeid", "hhid", "coupleid", "wave", "isced1997_r")
369 wave8gv_isced<-gv_isced[,selectr]
370
371 iv<-read_dta('sharew8_rel8-0-0_iv.dta')
372 names(iv)[names(iv)=="hhid8"]="hhid"
373 names(iv)[names(iv)=="mergeidp8"]="mergeidp"
374 names(iv)[names(iv)=="coupleid8"]="coupleid"
375 iv$wave<-8
376 selectiv<-c("mergeid", "hhid", "coupleid", "wave", "iv009_")
377 wave8iv<-iv[,selectiv]
378
379
380 im <- read_dta("sharew8_rel8-0-0_gv_imputations.dta")
381 im$wave=8
382 names(im)[names(im)=="hhid8"]="hhid"
383 names(im)[names(im)=="mergeidp8"]="mergeidp"
384 names(im)[names(im)=="coupleid8"]="coupleid"
385 colSums(is.na(im))
386 wave8im<-im[,c("mergeid", "hhid", "coupleid", "wave", "implicat", "thinc",
   sphus", "mstat", "nchild", "chronic",

```

```

387         "esmoked", "eurod", "doctor", "rhfo", "ghfo", "fdistress"
388         , "yedu" ) ]
389 wave8im$fdistress [ wave8im$fdistress == -99 ] = NA
390 wave8im$esmoked [ wave8im$esmoked == -99 ] = NA
391 wave8im$ghfo [ wave8im$ghfo == -99 ] = NA
392 wave8im = wave8im %>% group_by(mergeid) %>%
393   dplyr::summarise( thinc_m = mean( thinc ) ,
394                     nchild_m = as.integer( mean( nchild ) ) , chronic_m = as.
395                       integer( mean( chronic ) ) )
396                     doctor_m = as.integer( mean( doctor ) ) , rhfo_m = as.integer(
397                       mean( rhfo ) ) , ghfo_m = as.integer( mean( ghfo ) ) ,
398                     fdistress_m = as.integer( mean( fdistress ) ) , yedu_m = mean(
399                       yedu ) )
400 wave8im <- wave8im %>%
401   mutate(
402     sphus_m = as.factor( sphus_m ) ,
403     mstat_m = as.factor( mstat_m ) ,
404     esmoked_m = as.factor( esmoked_m ) ,
405     rhfo_m = as.factor( rhfo_m ) ,
406     ghfo_m = as.factor( ghfo_m ) ,
407     fdistress_m = as.factor( fdistress_m )
408   )
409 wave8im$wave = 8
410
411 m1 = merge( wave8co , wave8gv_iscd , by = intersect( names( wave8co ) , names(
412   wave8gv_iscd ) ) , all = TRUE )
413 m2 = merge( m1 , wave8iv , by = intersect( names( m1 ) , names( wave8iv ) ) , all = TRUE
414 )
415 wave84 = merge( m2 , wave8im , by = intersect( names( m2 ) , names( wave8im ) ) , all =
416   TRUE )
417
418 #write.csv( wave84 , file = "wave84.csv" , row.names = F )
419
420 #####
421 #####Merge
422 MM1 <- merge( wave81new , wave82 , by = intersect( names( wave81new ) , names(
423   wave82 ) ) , all = TRUE )
424 MM2 <- merge( MM1 , wave83new3 , by = intersect( names( MM1 ) , names( wave83new3 ) ) ,
425   all = TRUE )
426 wave8 <- merge( MM2 , wave84 , by = intersect( names( MM2 ) , names( wave84 ) ) , all =
427   TRUE )
428
429 #write.csv( wave8 , file = "wave8.csv" , row.names = F )
430 wave1_8 = merge( data , wave8 , by = intersect( names( data ) , names( wave8 ) ) , all =
431   TRUE )

```

B.3 Load data and data clearing

```

1 library( "RColorBrewer" )
2 library( "rattle" )
3 library( fdapace )

```

```

4 library("ggplot2")
5 library(caret)
6 library(tidyr)
7 library(glmnet)
8 library(viridis)
9 library(rpart)
10 library(rpart.plot)
11 library(randomForest)
12 library(cvTools)
13
14 #read data
15 data=read.csv("wave1_8_5.24.2022.csv")
16
17 #set the color
18 mycolor<-brewer.pal(9, "Set1")
19 mycolor.alpha <- scales::alpha(mycolor, 80/100)
20 mygray <- scales::alpha("gray", 40/100)
21 myeffectcolor<-viridis_pal(option = "C")(3)[c(1, 3, 2)]
22 myeffectcolor[2] <- "gray"
23 mycolor.alpha <- scales::alpha(mycolor, 80/100)
24
25
26 # set all the above three values to be NA.
27 RmID <- ((ID1 + ID2 + ID3) > 0)
28 myvar[RmID] = NA
29 return(myvar)
30
31
32 #Only keep subjects without missing cognitive functions
33 subdata <- subset(data, (recall_1>=0)&(recall_2>=0))
34
35 #Only keep the people age 50 and older
36 baseline <- subset(subdata, wave == 1)
37 ID1 <- baseline$mergeid[which(baseline$age >=50)]
38 subdata <- subset(subdata, mergeid%in%ID1)
39
40 #Scale my cognitive functions
41 subdata[, c("recall_1.scale", "recall_2.scale")] <- scale(subdata[, c("
  recall_1", "recall_2")])
42
43 # Only keep those with 7 measurements
44 myt <- table(subdata$mergeid)
45 mytt <- table(myt)
46 length(myt)
47 print(mytt/length(myt)*100)
48
49 myID <- names(myt)[myt >= 7]
50 subdata1 <- subset(subdata, mergeid%in%myID)
51 dim(subdata1)
52 #22876      54
53
54 #only keep the age bw. 60 and 85 in wave1
55 ID <- baseline$mergeid[which((baseline$age_int <= 80)&(baseline$age_int
  >= 65))]
```

```

56
57 #exclude the age over 95 in the last wave(the sample size is too small,
    outlier)
58 wave8 <- subset(subdata, wave == 8)
59 ID2 <- wave8$mergeid[which(wave8$age_int < 95)]
60
61 subdata2 <- subset(subdata1, mergeid%in%ID)
62 subdata2 <- subset(subdata2, mergeid%in%ID2)
63 dim(subdata2)
64 #6608    55
65
66 # baseline covariates
67 base02=subset(subdata2, select = -c(recall_1,recall_2))
68 base2 <- subset(base02, wave == 1)
69
70 colnames(base2)[2:53] <- paste("base.", colnames(base2)[2:53], sep = "")
71 }

```

B.4 FPCA

```

1 #FPCA
2 BFPCA2 <- MakeFPCAInputs(IDs = subdata2$mergeid, subdata2$age_int,
    subdata2$recall_2.scale)
3 FPCA2 <- FPCA(BFPCA2$Ly, BFPCA2$Lt)
4
5 plot(FPCA2)
6
7
8 #Figure 4.1 code
9 par(mfrow = c(1, 3))
10 #time #mean
11 plot(FPCA2$workGrid, FPCA2$mu, type = "l", xlab = "age",
12      ylab = "cognitive_function", lwd = 2, col = mycolor[1], main = "
    Mean_Function")
13 CreateScreePlot(FPCA2)
14 plot(FPCA2$workGrid, FPCA2$phi[,1], type = "l", col = mycolor[1],
15      ylim = c(-0.4, 0.4), lwd = 2, xlab = "age", ylab = "cognitive_
    function", lty = 1,
16      main = "First_3_Eigenfunctions")
17 lines(FPCA2$workGrid, FPCA2$phi[,2], col = mycolor[4], lwd = 2, lty = 2)
18 lines(FPCA2$workGrid, FPCA2$phi[,3], col = mycolor[8], lwd = 2, lty = 3)
19 abline(h = 0, col = "gray")
20 legend("top", paste(c("First", "Second", "Third"), "Eigenfunction"),
21       lwd=2, lty = 1:3, col = mycolor[c(1, 4, 8)], bty = "n")
22
23
24
25 #Figure effect of the First FPC
26 par(mfrow=c(1,2))
27 plot(FPCA2$workGrid, FPCA2$mu, type = "l", xlab = "age", ylab = "
    cognitive_function",
28      main = "Effect_of_First_Eigenfunction", lwd = 1.5)

```



```

29 lines(FPCA2$workGrid, FPCA2$mu - 0.3*FPCA2$phi[,1], col = myeffectcolor
    [1], lwd = 1.5)
30 lines(FPCA2$workGrid, FPCA2$mu + 0.3*FPCA2$phi[,1], col = myeffectcolor
    [3], lwd = 1.5)
31
32 plot(FPCA2$workGrid, FPCA2$mu, type = "l", xlab = "age", ylab = "
    cognitive_function",
33     main = "Effect_of_Second_Eigenfunction", lwd = 1.5)
34 lines(FPCA2$workGrid, FPCA2$mu - 0.3*FPCA2$phi[,2], col = myeffectcolor
    [1], lwd = 1.5)
35 lines(FPCA2$workGrid, FPCA2$mu + 0.3*FPCA2$phi[,2], col = myeffectcolor
    [3], lwd = 1.5)
36
37
38 #merge FPC1 and FPC2 with baseline covariates dataset
39 First2PC2 <- data.frame(cbind(names(BFPCA2$Ly), FPCA2$xiEst[, 1:2]))
40 colnames(First2PC2) <- c("mergeid", "PC1", "PC2")
41 First2PC2$PC1 <- as.numeric(as.character(First2PC2$PC1))
42 First2PC2$PC2 <- as.numeric(as.character(First2PC2$PC2))
43 PCs2 <- merge(First2PC2, base2, by = "mergeid")

```

B.5 PC score V.s. nine countries

```

1  ###PC1 for nine countries (Figure 4.4)
2  ggplot(PCs2, aes(x=country, y=PC1)) +
3    geom_boxplot(fill = mycolor, width = 0.6, notch = T)+
4    geom_hline(yintercept = 0, col = "gray") +
5    theme(axis.text.x = element_text(face = "bold", angle = 70, size = 8),
6          axis.ticks.x = element_blank(),
7          axis.text.y = element_text(face = "bold", size = 8),
8          panel.background = element_rect(fill = "white", colour = NA),
9          panel.border = element_rect(fill = "NA", color = "gray"),
10         panel.grid.major = element_line(colour = "gray80", size = 0.25),
11         panel.grid.minor = element_line(colour = "gray80", size = 0.25)
12    )+
13    labs(x="", y="First_PC_Score")
14
15  ###PC2 for nine countries (Figure 4.5)
16  ggplot(PCs2, aes(x=country, y=PC2)) +
17    geom_boxplot(fill = mycolor, width = 0.6, notch = T)+
18    geom_hline(yintercept = 0, color = "gray") +
19    theme(axis.text.x = element_text(face = "bold", angle = 70, size = 8),
20          axis.ticks.x = element_blank(),
21          axis.text.y = element_text(face = "bold", size = 8),
22          panel.background = element_rect(fill = "white", colour = NA),
23          panel.border = element_rect(fill = "NA", color = "gray"),
24          panel.grid.major = element_line(colour = "gray80", size = 0.25),
25          panel.grid.minor = element_line(colour = "gray80", size = 0.25)
26    )+labs(x="", y="Second_PC_Score")
27  Median11 <- with(PCs2, tapply(PC1, country, median))##median of PC1 for
    each country

```

```

28 Median22 <- with(PCs2, tapply(PC2, country, median))##median of PC2 for
    each country
29
30
31 mymedcurves22 <- data.frame(cbind(FPCA2$workGrid,FPCA2$mu+
32                                outer(FPCA2$phi[,1], Median11))) #
    eigenfunction one times median11
33
34
35 data_long22 <- gather(mymedcurves22, country, value, Austria:Switzerland
    , factor_key=TRUE)
36 ### Figure 4.2 code
37 ggplot(data_long22, aes(x=V1, y=value, group = country, color = country)
    )+
38   geom_line()+
39   scale_color_manual(values=mycolor) +
40   labs(x = "age", y = "cognitive_function")
41
42 mymedcurves33 <- data.frame(cbind(FPCA2$workGrid,
43                                FPCA2$mu+
44                                # outer(FPCA2$phi[,1], Median11)+
45                                outer(FPCA2$phi[,2], Median22)))
46 data_long33 <- gather(mymedcurves33, country, value, Austria:Switzerland
    , factor_key=TRUE)
47
48
49 ggplot(data_long33, aes(x=V1, y=value, group = country, color = country)
    )+
50   geom_line()+
51   scale_color_manual(values=mycolor) +
52   labs(x = "age", y = "cognitive_function")
53
54 ### Figure 4.3 code
55 data_long44 <- gather(mymedcurves44, country, value, Austria:Switzerland
    , factor_key=TRUE)
56
57 data_long44$country <- factor(data_long44$country, levels = c("Sweden",
    "Switzerland", "Denmark",
58                                                                "Germany",
    "France",
    "Austria",
    "Belgium",
    "Italy", "Spain"
59                                                                ))
60 pdf("median3.pdf")
61 ggplot(data_long44, aes(x=V1, y=value, group = country))+
62   geom_line(aes(linetype=country, color = country))+
63   scale_color_manual(values=c("#000000", "#E69F00", "#56B4E9",
64                                "#009E73", "red", "#0072B2", "#D55E00",

```

```

65         "#CC79A7", "#999999")))+
66     labs(x = "age", y = "cognitive_function")
67 dev.off()

```

B.6 data management before model fitting

```

1  #select covariates
2  select.co<-c("base.bmi", "base.thinc_m", "base.yedu_m", "base.maxgrip", "
    base.hhsz", "base.recall_1.scale", "base.recall_2.scale", 'base.
    siblings_alive')
3  select.dis <- c("base.country", "base.gender", "base.hearing", "base.
    reading",
4      "base.writing", "base.eyesight", "base.ep005", "base.mob
    ",
5  "base.active2", "base.active", "base.doctorS", "base.depression",
6  "base.chronic", "base.health", "base.ghelp", "base.rhelp",
7  "base.esmoked_m", "base.marry", "base.hospital_m",
8  "base.partnerinh", "base.mother_alive", "base.father_alive",
9  "base.hc029_", "base.grossmotor", "base.finemotor", "base.orienti",
10 'base.numeracy', "base.drink", "base.muscle", "base.vigorous")
11
12 mynewPCs0 <- PCs2
13 mynewPCs0[select.dis] <- lapply(mynewPCs0[select.dis], factor)
14
15 #new dataset with the select coulmn
16 mynewPCs3=cbind(mynewPCs0[1:3], mynewPCs0[, c(select.dis, select.co)])
17 mynewPCs3=mynewPCs3[complete.cases(mynewPCs3), ]

```

B.7 Model fitting

B.7.1 Elastic Net

```

1  # Fitted Model based on minimum Mean Squared Error
2  set.seed(1000)
3  myfit3 <- train(
4      PC1 ~.-mergeid-PC2, data = mynewPCs3, method = "glmnet",
5      trControl = trainControl("cv", number = 10),
6      tuneLength = 30
7  )
8
9  get_best_result(myfit3)
10 mycoef13 <- coef(myfit3$finalModel, myfit3$bestTune$lambda)
11
12 mycoef3 <- data.frame(rep(0, length(mycoef13@Dimnames[[1]])),
13     mycoef13@Dimnames[[1]])
14 mycoef3[mycoef13@i + 1, 1] <- mycoef13@x
15 colnames(mycoef3) <- c("V1", "V2")
16
17 newcoef3 <- mycoef3[order(mycoef3$V1), ]

```

```

17 newcoef3$V2 <- factor(newcoef3$V2, level = unique(as.character(newcoef3$
    V2)))
18 newcoef3$color <- myeffectcolor[sign(newcoef3$V1) + 2]
19
20 ## Figure 4.6 code
21 ggplot(newcoef3, aes(x=V2, y=V1, fill=color, color = color)) +
22   geom_bar(stat="identity")+
23   coord_flip()+
24   scale_x_discrete(position = "top")+
25   scale_fill_manual(name = "",
26                     labels = c("Negative_Effect", "Positive_Effect", "No
    _Effect"),
27                     values=myeffectcolor[c(1,3,2)])+
28   scale_color_manual(name = "",
29                      labels = c("Negative_Effect", "Positive_Effect", "
    No_Effect"),
30                      values=myeffectcolor[c(1,3,2)])+
31   theme(#axis.text.x = element_blank(),
32         axis.text.x = element_text(face = "bold"),
33         axis.text.y = element_text(face = "bold", color = newcoef3$color),
34         legend.position="bottom")+
35   labs(x = "", y = "", title = "Covariate_Effects_on_the_First_FPC_
    Scores")
36 #####
37 #####
38 set.seed(100000)
39 myfit23 <- train(
40   PC2 ~.-mergeid-PC1, data = mynewPCs3, method = "glmnet",
41   trControl = trainControl("cv", number = 10),
42   tuneLength = 10
43 )
44 mycoef23 <- coef(myfit23$finalModel, myfit23$bestTune$lambda)
45 mycoef2.23 <- data.frame(rep(0, length(mycoef23@Dimnames[[1]])),
    mycoef23@Dimnames[[1]])
46 mycoef2.23[mycoef23@i + 1, 1] <- mycoef23@x
47 colnames(mycoef2.23) <- c("V1", "V2")
48
49 newcoef23 <- mycoef2.23[order(mycoef2.23$V1), ]
50 newcoef23$V2 <- factor(newcoef23$V2, level = unique(as.character(
    newcoef23$V2)))
51 newcoef23$color <- myeffectcolor[sign(newcoef23$V1) + 2]
52 #newcoef2[which(sign(newcoef2$V1)==-1),]
53 #newcoef2[which(sign(newcoef2$V1)==1),]
54
55 ## Figure 4.7 code
56 ggplot(newcoef23, aes(x=V2, y=V1, fill=color, color = color)) +
57   geom_bar(stat="identity")+
58   coord_flip()+
59   scale_x_discrete(position = "top")+
60   scale_fill_manual(name = "",
61                     labels = c("Negative_Effect", "Positive_Effect", "No
    _Effect"),
62                     values=myeffectcolor[c(1,3,2)])+
63   scale_color_manual(name = "",

```

```

64         labels = c("Negative_Effect", "Positive_Effect", "
                    No_Effect"),
65         values=myeffectcolor[c(1,3,2)])+
66 theme(text = element_text(size=10),
67       axis.text.x = element_text(face = "bold"),
68       axis.text.y = element_text(face = "bold", color = newcoef23$color),
69       legend.position="bottom")+
70 labs(x = "", y = "", title = "Covariate_Effects_on_the_Second_FPC_
    Scores")

```

B.7.2 Decision Tree model with discrete variables

```

1  # fit a big tree
2  dat02 <- subset(mynewPCs3, select = c("PC1", select.dis))
3  dat02[2:31] <- lapply(dat02[2:31], factor)
4
5  fit02 <- rpart(PC1~., data=dat02)
6  model02 <- prune(fit02, cp=fit02$cptable[which.min(fit02$cptable[, "
    xerror"]),"CP"])
7  # This function returns the optimal cp value associated with the minimum
    error.
8  rpart.plot(model02, type=0,cex=1.2)
9
10 #Figure A.1
11 fancyRpartPlot(model02, uniform=TRUE,sub="")
12
13 dat12 <- subset(mynewPCs3, select = c("PC2", select.dis))
14 dat12[2:31] <- lapply(dat12[2:31], factor)
15
16 fit12 <- rpart(PC2~., data=dat12)
17 model12 <- prune(fit12, cp=fit12$cptable[which.min(fit12$cptable[, "
    xerror"]),"CP"])
18
19 rpart.plot(model12, type=0,cex=1.2)
20
21 #Figure A.2
22 fancyRpartPlot(model12, uniform=TRUE,sub="")

```

B.7.3 Decision Tree model with whole selected variables

```

1  dat002 <- subset(mynewPCs3, select = c("PC1", select.dis, select.co))
2  fit02 <- rpart(PC1~., data=dat002)
3  model002 <- prune(fit02, cp=fit02$cptable[which.min(fit02$cptable[, "
    xerror"]),"CP"])
4  rpart.plot(model002, type=0,cex=1.2)
5
6  #Figure 4.8
7  fancyRpartPlot(model002, uniform=TRUE,sub="")
8
9  dat112 <- subset(mynewPCs3, select = c("PC2", select.dis, select.co))
10 fit112 <- rpart(PC2~., data=dat112)

```

```

11 model112 <- prune(fit112 , cp=fit112$cptable[ which.min(fit112$cptable[, "
      xerror" ]), "CP" ])
12 rpart.plot(model112, type=0,cex=1.2)
13
14 #Figure 4.9
15 fancyRpartPlot(model112, uniform=TRUE, sub="")

```

B.7.4 Random Forest with whole selected variables

```

1 fit2 <- randomForest(PC1 ~ ., importance = TRUE, data=dat002)
2 print(fit2)
3 fit12 <- randomForest(PC2 ~ ., importance = TRUE, data=dat112)
4 print(fit12)
5 #Figure 4.10
6 par(mfrow = c(1,2))
7 varImpPlot(fit2, main = "For_the_First_FPC_Score", type=1)
8 varImpPlot(fit12, main = "For_the_Second_FPC_Score", type=1)

```

B.8 Predictive R^2 for different models

B.8.1 Elastic-net

```

1 #For FPC1
2 submynewPCs3=subset(mynewPCs3, select = -c(mergeid, PC2))
3 k <- 10 #the number of folds
4 folds <- cvFolds(NROW(submynewPCs3), K=k)
5 submynewPCs3$holdoutpred <- rep(0, nrow(submynewPCs3))
6
7 for(i in 1:k){
8   train <- submynewPCs3[ folds$subsets[ folds$which != i ], ] #Set the
      training set
9   validation <- submynewPCs3[ folds$subsets[ folds$which == i ], ] #Set the
      validation set
10   set.seed(1000000)
11   newlm <- train(
12     PC1 ~.-holdoutpred, data = train, method = "glmnet", #change to -
      holdoutpred-base.recall_2.scale
13     tuneGrid =expand.grid(alpha=0.1,lambda = 0.190567)) #the best tuning
      parameter
14   newpred <- predict(newlm, newdata=validation)
15   submynewPCs3[ folds$subsets[ folds$which == i ], ]$holdoutpred <- newpred
16 }
17
18 submynewPCs3$holdoutpred #do whatever you want with these predictions
19 ppress=sum((submynewPCs3$PC2 - submynewPCs3$holdoutpred)^2)
20 ss=sum((submynewPCs3$PC2-mean(submynewPCs3$PC2))^2)
21 rsquare=1-ppress/ss
22 rsquare
23 #0.45 elastic net predictive R^2 for PC2 without drs
24 #0.55 elastic net predictive R^2 for PC2 with drs
25 #####

```

```

26 #For FPC2
27 submynewPCs3=subset(mynewPCs3, select = -c(mergeid, PC1))
28 k <- 10 #the number of folds
29 folds <- cvFolds(NROW(submynewPCs3), K=k)
30 submynewPCs3$holdoutpred <- rep(0, nrow(submynewPCs3))
31
32 for(i in 1:k){
33   train <- submynewPCs3[ folds$subsets[ folds$which != i], ] #Set the
      training set
34   validation <- submynewPCs3[ folds$subsets[ folds$which == i], ] #Set the
      validation set
35   set.seed(1000000)
36   newlm <- train(
37     PC2 ~.-holdoutpred, data = train, method = "glmnet", #change to -
      holdoutpred-base.recall_2.scale
38     tuneGrid =expand.grid(alpha=0.1,lambda = 0.190567)) #the best tuning
      parameter
39   newpred <- predict(newlm, newdata=validation)
40   submynewPCs3[ folds$subsets[ folds$which == i], ]$holdoutpred <- newpred
41 }
42
43 submynewPCs3$holdoutpred #do whatever you want with these predictions
44 ppress=sum((submynewPCs3$PC2 - submynewPCs3$holdoutpred)^2)
45 ss=sum((submynewPCs3$PC2-mean(submynewPCs3$PC2))^2)
46 rsquare=1-ppress/ss
47 rsquare
48 #0.115 elastic net predictive R^2 for PC2 without drs
49 #0.132 elastic net predictive R^2 for PC2 with drs

```

B.8.2 Lasso

```

1 ##select best tuning parameter lambda
2 set.seed(10000)
3 lasso1 <- train(
4   PC1 ~.-mergeid-PC2-base.recall_2.scale, data = mynewPCs3, method = "
      glmnet",
5   trControl = trainControl(method = "repeatedcv",
6                             number = 10,
7                             repeats = 5,
8                             verboseIter = TRUE),
9   tuneGrid = expand.grid(alpha = 1, lambda=seq(0, 0.5, by = 0.1)),
10  tuneLength = 10)
11 get_best_result(lasso1)
12
13 datlas <- subset(mynewPCs3, select = c("PC1", select.dis, select.co))
14
15 #For PFC1
16 datlas <- subset(mynewPCs3, select = c("PC1", select.dis, select.co))
17
18 k <- 10 #the number of folds
19 folds <- cvFolds(NROW(datlas), K=k)
20 datlas$holdoutpred <- rep(0, nrow(datlas))
21

```

```

22 for(i in 1:k){
23   train <- datlas[folds$subsets[folds$which != i], ] #Set the training
      set
24   validation <- datlas[folds$subsets[folds$which == i], ] #Set the
      validation set
25   set.seed(10000)
26   newlm <- train(PC1 ~.-holdoutpred, data = train, method = "glmnet",
27     tuneGrid =expand.grid(alpha=1,lambda = 0.1)) #change to -holdoutpred
      -base.recall_2.scale
28   newpred <- predict(newlm,newdata=validation)
29   datlas[folds$subsets[folds$which == i], ]$holdoutpred <- newpred
30 }
31
32 datlas$holdoutpred #do whatever you want with these predictions
33 ppress=sum((datlas$PC1 - datlas$holdoutpred)^2)
34 ss=sum((datlas$PC1-mean(datlas$PC1))^2)
35 rsquare=1-ppress/ss
36 rsquare
37 ##0.54 lasso predictive R^2 for PC1 with drs
38 ##0.4 lasso predictive R^2 for PC1 without drs
39 #####
40 datlas2 <- subset(mynewPCs3, select = c("PC2", select.dis, select.co))
41
42 k <- 10 #the number of folds
43 folds <- cvFolds(NROW(datlas2), K=k)
44 datlas2$holdoutpred <- rep(0,nrow(datlas2))
45
46 for(i in 1:k){
47   train <- datlas2[folds$subsets[folds$which != i], ] #Set the training
      set
48   validation <- datlas2[folds$subsets[folds$which == i], ] #Set the
      validation set
49   set.seed(10000)
50   newlm <- train(PC2 ~.-holdoutpred, data = train, method = "glmnet",
51     tuneGrid =expand.grid(alpha=1,lambda = 0.1)) #change to
      -holdoutpred-base.recall_2.scale
52   newpred <- predict(newlm,newdata=validation)
53   datlas2[folds$subsets[folds$which == i], ]$holdoutpred <- newpred t
54 }
55
56 datlas2$holdoutpred #do whatever you want with these predictions
57 ppress=sum((datlas2$PC2 - datlas2$holdoutpred)^2)
58 ss=sum((datlas2$PC2-mean(datlas2$PC2))^2)
59 rsquare=1-ppress/ss
60 rsquare
61 ##0.13 lasso predictive R^2 for PC2 with drs
62 ##0.08 lasso predictive R^2 for PC2 without drs

```

B.8.3 Random Forest

```

1 datc00 <- subset(mynewPCs3, select = c("PC1", select.dis, select.co))
2
3 k <- 10 #the number of folds

```



```

4 folds <- cvFolds(NROW(datc00), K=k)
5 datc00$holdoutpred <- rep(0,nrow(datc00))
6
7 for(i in 1:k){
8   train <- datc00[folds$subsets[folds$which != i], ] #Set the training
   set
9   validation <- datc00[folds$subsets[folds$which == i], ] #Set the
   validation set
10  set.seed(10000)
11  newlm <- randomForest(PC1 ~ .-holdoutpred, data=train) #change to -
   holdoutpred-base.recall_2.scale
12  newpred <- predict(newlm, newdata=validation)
13  datc00[folds$subsets[folds$which == i], ]$holdoutpred <- newpred
14 }
15 datc00$holdoutpred #do whatever you want with these predictions
16 ppress=sum((datc00$PC1 - datc00$holdoutpred)^2)
17 ss=sum((datc00$PC1-mean(datc00$PC1))^2)
18 rsquare=1-ppress/ss
19 rsquare
20 #0.479 random forest predictive R^2 for PC1 with drs
21 #0.37 random forest predictive R^2 for PC1 without drs
22 #####
23 datc112 <- subset(mynewPCs3, select = c("PC2", select.dis, select.co))
24 folds <- cvFolds(NROW(datc112), K=k)
25 datc112$holdoutpred <- rep(0,nrow(datc112))
26
27 for(i in 1:k){
28   train <- datc112[folds$subsets[folds$which != i], ] #Set the training
   set
29   validation <- datc112[folds$subsets[folds$which == i], ] #Set the
   validation set
30   set.seed(1000000)
31   newlm <- randomForest(PC2 ~ .-holdoutpred, data=train) #change to -
   holdoutpred-base.recall_2.scale
32   newpred <- predict(newlm, newdata=validation)
33   datc112[folds$subsets[folds$which == i], ]$holdoutpred <- newpred }
34
35 datc112$holdoutpred #do whatever you want with these predictions
36 ppress=sum((datc112$PC2 - datc112$holdoutpred)^2)
37 ss=sum((datc112$PC2-mean(datc112$PC2))^2)
38 rsquare=1-ppress/ss
39 rsquare
40 #0.122 random forest predictive R^2 for PC2 with drs
41 #0.1 random forest predictive R^2 for PC2 without drs

```

B.8.4 Decision Tree

```

1 k <- 10 #the number of folds
2 folds <- cvFolds(NROW(datc00), K=k)
3 datc00$holdoutpred <- rep(0,nrow(datc00))
4
5 for(i in 1:k){

```

```

6   train <- datc00[folds$subsets[folds$which != i], ] #Set the training
    set
7   validation <- datc00[folds$subsets[folds$which == i], ] #Set the
    validation set
8   set.seed(10000)
9   newlm <- prune(rpart(PC1~.-holdoutpred, data=train), cp=rpart(PC1~.-
    holdoutpred, data=train)
10  $cptable[which.min(rpart(PC1~.-holdoutpred, data=train)$cptable[, "
    xerror"]),"CP"])
11  #change to -holdoutpred-base.recall_2.scale
12  newpred <- predict(newlm,newdata=validation)
13  datc00[folds$subsets[folds$which == i], ]$holdoutpred <- newpred
14  }
15
16  datc00$holdoutpred #do whatever you want with these predictions
17  ppress=sum((datc00$PC1 - datc00$holdoutpred)^2)
18  ss=sum((datc00$PC1-mean(datc00$PC1))^2)
19  rsquare=1-ppress/ss
20  rsquare
21  #0.436 decision tree predictive R^2 for PC1 with drs
22  #0.26 decision predictive R^2 for PC1 without drs
23  #####
24  for(i in 1:k){
25    train <- datc112[folds$subsets[folds$which != i], ] #Set the training
        set
26    validation <- datc112[folds$subsets[folds$which == i], ] #Set the
        validation set
27    set.seed(1000000)
28    newlm <- prune(rpart(PC2~.-holdoutpred, data=train), cp=rpart(PC2~.-
        holdoutpred, data=train)$cptable[which.min(rpart(PC2~.-holdoutpred
        , data=train)$cptable[, "xerror"]),"CP"])
29    #change to -holdoutpred-base.recall_2.scale
30    newpred <- predict(newlm,newdata=validation)
31    datc112[folds$subsets[folds$which == i], ]$holdoutpred <- newpred
32  }
33
34  datc112$holdoutpred #do whatever you want with these predictions
35  ppress=sum((datc112$PC2 - datc112$holdoutpred)^2)
36  ss=sum((datc112$PC2-mean(datc112$PC2))^2)
37  rsquare=1-ppress/ss
38  rsquare
39  #0.094 decision tree predictive R^2 for PC2 with drs
40  #0.074 decision tree predictive R^2 for PC2 without drs

```