



Machine Learning Workflow to Explain Black-Box Models for Early Alzheimer's Disease Classification Evaluated for Multiple Datasets

Louise Bloch^{1,2} · Christoph M. Friedrich^{1,2} · for the Alzheimer's Disease Neuroimaging Initiative

Received: 11 October 2021 / Accepted: 14 August 2022 / Published online: 6 October 2022
© The Author(s) 2022

Abstract

Hard-to-interpret black-box Machine Learning (ML) was often used for early Alzheimer's Disease (AD) detection. To interpret eXtreme Gradient Boosting (XGBoost), Random Forest (RF), and Support Vector Machine (SVM) black-box models, a workflow based on Shapley values was developed. All models were trained on the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset and evaluated for an independent ADNI test set, as well as the external Australian Imaging and Lifestyle flagship study of Ageing (AIBL), and Open Access Series of Imaging Studies (OASIS) datasets. Shapley values were compared to intuitively interpretable Decision Trees (DTs), and Logistic Regression (LR), as well as natural and permutation feature importances. To avoid the reduction of the explanation validity caused by correlated features, forward selection and aspect consolidation were implemented. Some black-box models outperformed DTs and LR. The forward-selected features correspond to brain areas previously associated with AD. Shapley values identified biologically plausible associations with moderate-to-strong correlations with feature importances. The most important RF features to predict AD conversion were the volume of the amygdalae and a cognitive test score. Good cognitive test performances and large brain volumes decreased the AD risk. The models trained using cognitive test scores significantly outperformed brain volumetric models ($p < 0.05$). Cognitive Normal (CN) vs. AD models were successfully transferred to external datasets. In comparison to previous work, improved performances for ADNI and AIBL were achieved for CN vs. Mild Cognitive Impairment (MCI) classification using brain volumes. The Shapley values and the feature importances showed moderate-to-strong correlations.

Keywords Interpretable machine learning · Early Alzheimer's disease detection · Shapley values

This article is part of the topical collection "Biomedical Engineering Systems and Technologies" guest edited by Hugo Gamboa and Ana Fred.

Membership of the Alzheimer's Disease Neuroimaging Initiative is listed in the Acknowledgments.

✉ Christoph M. Friedrich
christoph.friedrich@fh-dortmund.de

Louise Bloch
louise.bloch@fh-dortmund.de

¹ Department of Computer Science, University of Applied Sciences and Arts Dortmund, Emil-Figge-Str. 42, Dortmund 44227, Germany

² Institute for Medical Informatics, Biometry and Epidemiology (IMIBE), University Hospital Essen, Hufelandstr. 55, Essen 45147, Germany

Introduction

Alzheimer's Disease (AD) is a neurodegenerative disease [1] and the most frequent cause of dementia. As the number of dementia patients increases continuously, AD is a globally growing health problem [2]. Currently, there is no causal therapy to cure AD [1]. To recruit and monitor subjects for therapy studies, it is important to identify patients at risk to develop AD early and to develop preclinical markers. Subjects with cognitive impairments that do not interfere with everyday activities are considered as having Mild Cognitive Impairment (MCI) due to AD [3]. The risk to develop AD is increased for subjects with MCI in comparison to cognitively normal controls (CN). However, not all subjects with MCI prospectively convert to AD. One possibility for early AD detection is to find patterns distinguishing between progressive MCI subjects (pMCI) who will develop AD and subjects with stable MCI (sMCI).

Multiple Machine Learning (ML) workflows were implemented for this differentiation. Some used models like Decision Trees (DTs) or Logistic Regression (LR), which were interpretable by design. However, black-box models like eXtreme Gradient Boosting (XGBoost) [4], Random Forests (RFs) [5], or Convolutional Neural Networks (CNNs) [6] often outperform those models. Black-box models are designed to identify highly complex associations and are challenging to interpret. Thus, the risk of learning spurious decision functions caused by patterns occurring in the training dataset is increased for black-box models [7].

This research is an extended version of earlier work [8] and thus expands the previously developed ML workflow. The previously developed workflow enabled the interpretation of black-box models based on model-agnostic Shapley values. Shapley values give individual explanations for the prediction of each subject and visualize complex relationships between features and model predictions. In this research, the previous experiments are expanded using three AD datasets and three adjusted feature sets. In addition to the previously trained tree-based models, Support Vector Machines (SVMs) [9] and LR models were implemented and explained. In this work, Shapley-based explanations were compared to classical feature importance methods, absolute log odd's ratios, and permutation importance.

In comparison to previous work [8], an improvement of the classification results for ADNI and AIBL was achieved for the differentiation between Cognitive Normal (CN) controls and MCI subjects as well as for MCI vs. AD classification and models trained without cognitive test scores and validated for AIBL. Additionally, the ADNI and AIBL results achieved for sMCI vs. pMCI classification, trained with cognitive test scores, outperformed previous work.

This article is structured as follows: In “[Related Work](#)”, related work is described. Section “[Materials and Methods](#)” introduces the datasets and methods used to implement the ML workflow and the details of the experiments. Section “[Results](#)” elaborates on the experimental results. Those results are discussed including the mentioning of limitations in “[Discussion](#)”. Finally, “[Conclusion](#)” concludes the overall work.

Related Work

Interpretable ML was developed to explain black-box models [10]. As the heterogeneous etiology of AD is not completely understood yet, interpretability is important and enables the validation of the biological plausibility of ML models. Recently, some studies have used interpretable ML in AD detection.

For example, Long Short-Term Memory- (LSTM-) [11] based Recurrent Neural Networks (RNN) [12] were trained

to classify CN vs. MCI subjects in [13]. The experiments included multiple techniques to fuse socio-demographic and genetic data with Magnetic Resonance Imaging (MRI) scans. The resulting models were evaluated for two AD datasets—the AD subset [14] of the Heinz Nixdorf Risk Factors Evaluation of Coronary Calcification and Lifestyle (RECALL) (HNR) [15] (61 MCI and 59 CN) and 624 subjects (397 MCI, 227 CN) of the Alzheimer's Disease Neuroimaging Initiative (ADNI) [16] study phase 1. To visually explain individual model decisions, Gradient-weighted Class Activation Mapping (Grad-CAM) [17] was used. A focus on biologically plausible regions was observed.

Four heatmap visualization methods—sensitivity analysis [18], guided backpropagation [19], occlusion [20], and brain area occlusion inspired by [21]—were compared for 3D-CNNs in [22]. The CNN models were trained using 969 MRI scans of 344 ADNI subjects (151 CN, 193 AD). However, it was unclear whether the described workflow ensured independent training and test sets using multiple scans per subject [23]. Thus, the Cross-Validation (CV) accuracy of $77\% \pm 6\%$ might be affected by data leakage. All heatmaps focused on AD-related anatomical brain areas.

An interpretable deep learning model consisting of a Generative Adversarial Network [24] to extend the training dataset, a regression network to generate feature vectors from adjacent visits, and a classification model was introduced in [25]. First, the regression model iteratively estimated the feature vector at the following visit. The resulting feature vector was used as input for the classification model, which predicted the final diagnosis. To classify 101 pMCI vs. 115 sMCI ADNI subjects, longitudinal volumetric MRI features were used. The model outperformed SVMs and artificial neural networks.

A new interpretable model based on distinct weighted rules was introduced in [26] and evaluated for 151 subjects (97 AD and 54 CN) of the ADNI cohort. The framework is called Sparse High-order Interaction Model with Rejection option (SHIMR) and consists of two hierarchical stages. In the first stage, the interpretable model was trained using plasma features. The data of subjects with an unclear prediction in this stage were propagated to the second stage. In this stage, an SVM [9] was trained using invasive Cerebrospinal Fluid (CSF) markers. The evaluation included both CV and an independent test set. The described model reached an Area Under the Receiver-Operating characteristics Curve (AUROC) of 0.81 for the test set.

SHapley Additive exPlanations (SHAP) [27] were used in [28] to explain differences in models trained using coresets selection methods. The idea was to determine coresets of subjects with the most informative data. RF and XGBoost models were trained on these coresets to avoid overfitting and improve ML models. The results of Data Shapley [29] coresets selection were compared to Leave-One-Out [30] selection and random exclusion. All models were trained

Table 1 Summary of the related work

Ref.	Task	Subjects	Modality	ML method	Explanability method
[13]	CN vs. MCI	HNR: 61 MCI, 59 CN; ADNI-1: 397 MCI, 227 CN	MRI, socio-demography, ApoE	LSTM based RNN	GradCAM
[22]	CN vs. AD	ADNI: 151 CN, 193 AD	MRI	CNN	Sensitivity analysis, guided backpropagation, occlusion, brain area occlusion
[25]	sMCI vs. pMCI	ADNI: 101 pMCI, 115 sMCI	MRI volumes	Neural network	Intrinsic
[26]	CN vs. AD	ADNI: 54 CN, 97 AD	CSF, Plasma	SHIMR	Intrinsic
[28]	sMCI vs. pMCI	ADNI: 400 sMCI, 319 pMCI; AIBL: 16 sMCI, 12 pMCI	MRI volumes, demography, ApoE	RF, XGBoost	SHAP
[31]	CN vs. MCI, CN vs. AD, MCI vs. AD	ADNI: 148 CN, 147 MCI, 110 AD	Amyloid-PET, MRI, FDG- PET, CSF	RF, GTB	RF-Feature importance, SHAP
[33]	high vs. low risk	SHARE: 80,699 CN, 4,157 AD; PREVENT: 364 low risk, 109 high risk	socio-demography, lifestyle	RF, XGBoost	SHAP
[37]	CN vs. MCI vs. AD, sMCI vs. pMCI	ADNI: 294 CN, 254 sMCI, 232 pMCI, 268 AD	MRI, CSF, PET, cognitive tests, medical history, genetics	RF	Ensemble of surrogat models, SHAP

and validated for the ADNI dataset (400 sMCI, 319 pMCI) and externally validated for a subset of the AIBL dataset (16 sMCI, 12 pMCI). SHAP summary plots showed that models trained for both the entire training set and the coreset learned biologically plausible associations.

To examine the predictive influence of β -amyloid plaques, tau tangles, and neurodegeneration during the disease progression, RF feature importance was used in [31]. The experimental data included 405 ADNI subjects (148 CN, 147 MCI, 110 AD). β -amyloid Positron Emission Tomography (PET) detected β -amyloid plaques, invasive CSF features surrogated tau tangles, and MRI and Fluorodeoxyglucose (FDG) PET scans were used to determine neurodegeneration. The experimental results showed that models trained to classify the early AD stages preferred features representing tau tangles and β -amyloid plaques. Models trained to predict later stages favored surrogates for neurodegeneration. SHAP [27] and Gradient Tree Boosting (GTB) [32] reproduced those observations. The RF and the entire feature set reached accuracies of 73.17 % (CN vs. MCI), 71.01 % (MCI vs. AD), and 90.34 % (CN vs. AD).

SHAP values were also used in [33] to explain population-based and individual predictions of XGBoost models and RFs. Models were trained using socio-demographic and lifestyle factors to predict the patient's risk to develop AD based on medical history. Transfer learning applied information extracted from the Survey of Health, Ageing, and Retirement in Europe (SHARE) [34] (80,699 CN, 4,157 AD) to the PREVENT cohort [35] (109 subjects with high risk to develop AD, 364 subjects with low risk). The PREVENT cohort was

younger than the SHARE cohort. The models support the hypothesis that age is the most important risk factor in AD detection. Consistent with previous research [36], among other factors, less education, physical inactivity, diabetes, and infrequent social contact were identified as potential risk factors.

A two stage-based classification workflow that used SHAP values to interpret RFs was developed in [37]. In the first stage, CN vs. MCI vs. AD classification was performed. The second stage implemented the differentiation of sMCI and pMCI subjects. The models were based on multiple modalities including MRI, PET, CSF biomarkers, cognitive tests, medical history, genetics, and many more. The RFs were trained and tested using 1,048 subjects (294 CN, 254 sMCI, 232 pMCI, and 268 AD) of the ADNI dataset. For CN vs. MCI vs. AD classification, the model almost exclusively selected cognitive test scores as the most important features. The model learned bad cognitive test results increased the risk of AD and MCI. The most important features for sMCI vs. pMCI classification also were cognitive test scores followed by PET and MRI features. Bad cognitive test scores, small MRI volumes, and small PET uptakes were associated with disease progression (Table 1).

Materials and Methods

The ML workflow, implemented using the programming language Python v3.6.9 [38], is shown in Fig. 1. It enables the interpretation of black-box models trained to detect early

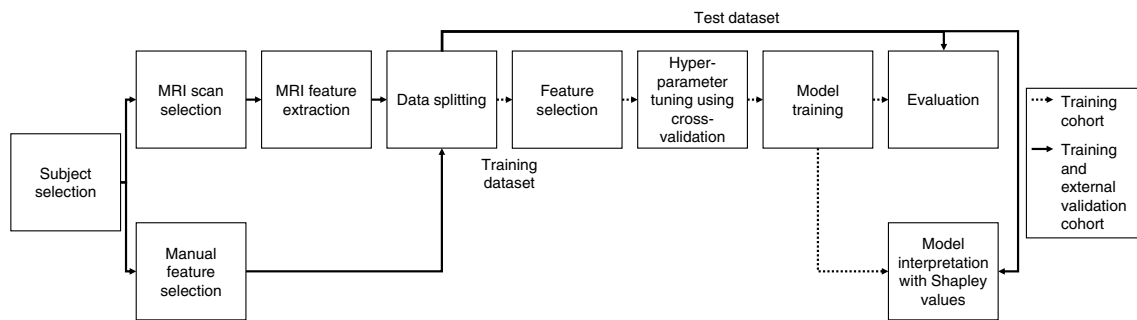


Fig. 1 Implemented ML workflow. Volumetric features were extracted for one baseline (BL) MRI scan per subject. The ADNI dataset was randomly split into an 80 % training and 20 % test set. The most important MRI features were selected using forward feature selection, and those were concatenated with socio-demographic features, number of ApolipoproteinE4 (ApoE4) alleles, and cognitive

AD. In the following, the workflow and the methods used for implementation are elucidated.

Datasets

Data used in the preparation of this article were obtained from the ADNI [16], the AIBL [39], and the OASIS [40] cohorts.

ADNI¹ was launched in 2003 as a public-private partnership. The primary goal of ADNI is to test whether a combination of biomarkers can measure the progression of MCI and AD. Those biomarkers include serial MRI, PET, biological markers, as well as clinical and neuropsychological assessments. The ongoing ADNI cohort recruited subjects from more than 60 sites in the United States and Canada and consists of four phases (ADNI-1, ADNI-2, ADNIGO, and ADNI-3). The subjects were assigned to three diagnostic groups. CNs have no problems with memory loss. Subjects with AD meet the criteria for probable AD defined by the National Institute of Neurological and Communicative Disorders and Stroke–Alzheimer’s Disease and Related Disorders Association (NINCDS-ADRDA) [41]. The diagnostic criteria of ADNI were explained in [16]. The dataset was downloaded on 27 Jul 2020 and initially included 2,250 subjects.

AIBL² is the largest AD study in Australia and was launched in 2006. AIBL aims to discover biomarkers, cognitive test results, and lifestyle factors associated with AD. As AIBL focuses on early AD stages, most of the subjects are CN. The MCI subjects of AIBL met the criteria described in [42], AD diagnoses following the NINCDS-ADRDA criteria [41] for probable AD. The diagnostic criteria of AIBL

test scores. Bayesian optimization implemented hyperparameter-tuning. Black-box RFs, XGBoost models, LR models, as well as polynomial and radial SVMs, were trained and validated. Shapley values were calculated for black-box model interpretation. An evaluation was performed for the independent ADNI test set and for the external AIBL and OASIS datasets

were described in [39]. Approximately half of the CN subjects recruited in AIBL show memory complaints [39]. AIBL data version 3.3.0 was downloaded on 19 Sep 2019 and originally included 858 subjects.

The aim of the Open Access Series of Imaging Studies (OASIS) 3³ [40] dataset is to investigate the effects of healthy ageing and AD. The subjects of OASIS-3 were recruited from several ongoing studies in the Washington University Knight Alzheimer Disease Research Center⁴. The longitudinal dataset included MRI scans, fMRI scans, Amyloid- and FDG-PET scans, neuropsychological test results, and clinical data for 1,098 subjects. OASIS focuses on the preclinical stage of AD. All OASIS subjects had a Clinical Dementia Rating (CDR) less than or equal to 1. The OASIS dataset provides multiple target values. In this research, CN subjects had normal cognition and absence of MCI or AD diagnosis, MCI subjects had amnesic MCI with memory impairment, and AD diagnosis follows the NINCDS-ADRDA criteria [41] for probable AD.

Subject Selection

For the ADNI dataset, all subjects with an MRI scan at the baseline visit were included. 521 subjects who have no MRI scan at the baseline visit were excluded, 29 subjects failed the MRI feature extraction described in “MRI Feature Extraction”. The demographics of the resulting 1,700 subjects are summarized in Table 2.

The 853 MCI subjects were divided into two groups. The sMCI subjects had a stable MCI diagnosis at all follow-up visits and the pMCI subjects converted to a stable AD

¹ ADNI: <https://adni.loni.usc.edu>, Accessed: 2022-05-01.

² AIBL: <https://aibl.csiro.au/>, Accessed: 2022-05-01.

³ OASIS: <https://www.oasis-brains.org/>, Accessed: 2022-05-01.

⁴ Washington University Knight Alzheimer Disease Research Center: <https://knightadrc.wustl.edu/>, Accessed: 2022-05-01.

Table 2 ADNI demographics at BL. The mean (\bar{x}) and standard deviation (σ) are given for all continuous variables

	<i>n</i>	Age years	Gender f in %	Education years	MMSCORE $\bar{x} \pm \sigma$	CDR $\bar{x} \pm \sigma$	ApoE ϵ 4 ¹ 0/1/2 in %
CN	512	74.2 \pm 5.8	51.8	16.3 \pm 2.7	29.1 \pm 1.1	0.0 \pm 0.0	71.3/26.2/ 2.3
MCI	853	73.1 \pm 7.6	40.8	15.9 \pm 2.9	27.6 \pm 1.8	0.5 \pm 0.0	49.4/39.5/10.8
sMCI	400	73.2 \pm 7.5	40.2	15.8 \pm 3.0	27.8 \pm 1.8	0.5 \pm 0.0	56.8/34.0/ 9.2
pMCI	319	74.0 \pm 7.1	40.1	15.9 \pm 2.8	27.0 \pm 1.7	0.5 \pm 0.0	34.2/49.5/16.3
AD	335	75.0 \pm 7.8	44.8	15.2 \pm 3.0	23.2 \pm 2.1	0.8 \pm 0.3	33.1/47.2/19.1
$\Sigma_{CN,MCI,AD}$	1,700	73.8 \pm 7.2	44.9	15.9 \pm 2.9	27.2 \pm 2.7	0.4 \pm 0.3	52.8/37.0/ 9.9

¹ For 6 ADNI subjects (1 CN, 3 MCI, 2 AD), the number of ApoE ϵ 4 alleles was missing

Table 3 AIBL demographics at BL. The mean (\bar{x}) and standard deviation (σ) are given for all continuous variables

	<i>n</i>	Age years	Gender f in %	MMSCORE $\bar{x} \pm \sigma$	CDR $\bar{x} \pm \sigma$	ApoE ϵ 4 ¹ 0/1/2 in %
CN	446	72.5 \pm 6.1	57.0	28.7 \pm 1.2	0.0 \pm 0.1	69.3/26.5/ 2.7
MCI	95	75.4 \pm 7.0	47.4	27.1 \pm 2.2	0.5 \pm 0.1	47.4/36.8/12.6
sMCI	16	77.8 \pm 6.9	37.5	28.0 \pm 1.7	0.4 \pm 0.2	56.2/37.5/ 6.2
pMCI	12	75.3 \pm 5.8	33.3	26.2 \pm 1.6	0.5 \pm 0.0	16.7/50.0/33.3
AD	71	73.1 \pm 6.6	59.2	20.5 \pm 5.7	0.9 \pm 0.6	29.6/49.3/18.3
$\Sigma_{CN,MCI,AD}$	612	73.0 \pm 6.6	55.7	27.5 \pm 3.5	0.2 \pm 0.4	61.3/30.7/ 6.0

¹ For 12 AIBL subjects (7 CN, 3 MCI, 2 AD), the number of ApoE ϵ 4 alleles was missing

Table 4 OASIS-3 demographics at the first visit with MRI scan and diagnosis

	<i>n</i>	Age years	Gender ¹ f in %	MMSCORE $\bar{x} \pm \sigma$	CDR $\bar{x} \pm \sigma$	ApoE ϵ 4 ² 0/1/2 %
CN	704	68.3 \pm 9.3	58.7	29.1 \pm 1.2	0.0 \pm 0.1	65.8/29.7/4.1
MCI	19	76.7 \pm 7.0	36.8	28.1 \pm 1.4	0.3 \pm 0.2	57.9/42.1/0.0
AD	198	75.6 \pm 7.9	48.5	24.8 \pm 4.0	0.7 \pm 0.3	38.9/51.5/9.1
$\Sigma_{CN,MCI,AD}$	921	70.1 \pm 9.5	56.0	28.1 \pm 2.8	0.2 \pm 0.3	59.8/34.6/5.1

The mean (\bar{x}) and standard deviation (σ) are given for all continuous variables

¹ For 1 OASIS subjects (1 CN), the gender was missing

² For 4 OASIS subjects (3 CN, 1 AD), the number of ApoE ϵ 4 alleles was missing

diagnosis at any visit. 38 subjects with no follow-up visits and 96 subjects who reverted to CN or MCI were excluded from this separation, resulting in 400 sMCI and 319 pMCI subjects.

For AIBL, the same exclusion criteria were applied. Therefore, 170 subjects had no MRI scan at the baseline visit, and the baseline MRI scans of 76 subjects failed for the MRI feature extraction pipeline described in “MRI Feature Extraction”. The demographics of the resulting 612 subjects are summarized in Table 3. Similar to the ADNI dataset, the 95 MCI subjects were divided into two groups. In this step, 60 subjects with no follow-up visits and 7 subjects who reverted to CN or MCI were excluded from this separation, resulting in 16 sMCI and 12 pMCI subjects.

The exclusion criteria were similarly applied for the OASIS-3 dataset, which originally included 1,098 subjects. For 983 subjects, a diagnosis of CN, MCI, or AD was assigned for at least one visit. The MRI feature extraction pipeline failed for all MRI scans of five subjects, and no MRI scan was successfully matched to a diagnosis with a tolerance of 365 days for 57 subjects. In contrast to the ADNI and AIBL datasets, which exclusively included baseline visits, the first visit with an MRI scan and a diagnosis was used for OASIS. The demographics of the remaining 921 subjects are summarized in Table 4.

The number of subjects with MCI as baseline diagnosis is 19. This number was decreased if subjects without follow-up diagnoses were excluded. Thus, no experiments were executed to separate sMCI and pMCI subjects in

OASIS-3. For reproducible research, the supplementary material contains lists with the subject and MRI IDs and the diagnoses for all datasets.

MRI Scan Selection

From the ADNI dataset, T1-weighted MRI scans recorded at the baseline visit were included. The acquisition parameters differ between scanners. During the ADNI-1 study phase, scans were recorded using a field strength of 1.5 T. In the remaining study phases, MRI scans with a field strength of 3.0 T were recorded.

From the AIBL dataset, T1-weighted MRI scans following the protocol of the ADNI 3D T1-weighted sequences were included. All AIBL scans had a resolution of $1 \times 1 \times 1.2$ mm.

For the OASIS-3 dataset, T1-weighted MRI scans, recorded on three scanners, were included. The field strengths of those scanners are 1.5 T and 3.0 T [40].

MRI Feature Extraction

Using FreeSurfer v6.0 [43], volumetric features were extracted for each MRI scan. These include the volumes of 34 cortical areas per hemisphere of the Desikan–Killiany atlas [44], 34 subcortical areas [45], and the estimated Total Intracranial Volume (eTIV). As recommended for volumes in [46], the volumetric features were normalized by eTIV. This results in 103 MRI volumes, which were split into 49 features of the left hemisphere, 49 features of the right hemisphere, and five additional unpaired segmentations (3rd ventricle, 4th ventricle, brain stem, CSF, eTIV).

After the normalization, for paired volumes, the sum (described in Eq. 1), the difference (described in Eq. 2), and the ratio (described in Eq. 3) of both hemispheres are calculated to investigate symmetry and to decrease feature interactions. This results in 152 MRI features (49 sums, 49 differences, 49 ratios, and 5 unpaired features). Brain asymmetry was previously associated with AD [47–50]. Equation 2 shows that differences were calculated by subtracting the right from the left volume similar to [48], where the cortical thickness was used instead of volumetric features

$$sum_{ROI} = \frac{vol_{ROI}^{left}}{eTIV} + \frac{vol_{ROI}^{right}}{eTIV} \quad (1)$$

$$diff_{ROI} = \frac{vol_{ROI}^{left}}{eTIV} - \frac{vol_{ROI}^{right}}{eTIV} \quad (2)$$

$$ratio_{ROI} = \frac{vol_{ROI}^{left}}{vol_{ROI}^{right}} \quad (3)$$

Manual Feature Preselection

Three feature sets were investigated in the experiments. The manual feature selection aims to choose less-invasive, accessible examination techniques which were able to detect early signs of AD. Feature set 1 (FS-1) includes all MRI features, and socio-demographic features including age, gender, and years of education. However, the years of education are only available for the ADNI dataset. Feature set 2 (FS-2) expands FS-1 by the number of ApoEε4 alleles, a genetic risk factor associated with AD, which can be obtained from blood samples or via less-invasive swab tests from the inside surface of the cheek. Feature set 3 (FS-3) extended FS-2 by three cognitive tests including the score of the Mini-Mental State Examination (MMSCORE) and two logical tests to evaluate the long-term (Logical memory, delayed—LDELTOTAL) and the short-term memory (Logical memory, immediate—LIMMTOTAL). The CDR was strongly associated with AD diagnosis and was not included in the experiments.

Dataset Splitting

The ADNI dataset was split on the subject level into two distinct subsets. The training set included 80 % of the data and the test set consisted of the remaining 20 %. The splitting was executed within each diagnostic group to ensure similar distributions. The AIBL and OASIS datasets were used as external test sets. None of the AIBL and OASIS subjects was used in the training or model selection process.

Feature Selection

Initially, 152 MRI features were extracted from the MRI scans. Those features are reduced to focus the ML models on the most important features. For this reason, feature forward selection was implemented. In comparison to feature selection methods like RF feature importance, this method avoids correlated features in the dataset [51]. Forward selection is a greedy procedure that iteratively identifies the best new feature until no improvement was reached. The training dataset was split into an 80 % training set and 20 % validation dataset. The training dataset was used to train the ML model used for classification with default hyperparameters on the feature set, and the validation dataset was used to calculate the validation accuracy for this feature set. The selected MRI features were expanded using the features described in “Manual Feature Preselection”.

Hyperparameter-Tuning

To tune the hyperparameters of the ML models, Bayesian optimization [52] was implemented using the Python package

scikit-optimize v0.8.1 [53]. Bayesian optimization maps the dependency of the hyperparameters and the model performance using a Gaussian Process. Initially, ten nearly random hyperparameter combinations were selected by a Latin Hypercube Design (LHD) [54]. Bayesian optimization with LHD initialization was successfully used in previous research [55] to optimize the parameters for early AD detection. Each parameter was split into ten equidistant intervals and one sample was randomly chosen per interval. This results in ten samples per parameter, which were randomly matched.

A stratified 10×10 -fold CV [56] was applied to the training dataset to estimate the model accuracy for an independent test set. Stratified 10×10 -fold CV was implemented by splitting each diagnostic group of the training dataset into ten distinct folds using the Python package scikit-learn v0.23.2 [57]. Ten iterations were performed, each with a different fold used as a validation dataset (10 %). The training dataset included the remaining 9 folds (90 %). With shuffled data in each run, this procedure was repeated ten times. The ML model was initially evaluated for ten LHD combinations.

To predict the average CV accuracy for the initial parameter combinations, the Gaussian process was fitted. Afterward, an optimization selected the next promising parameter combination. As an acquisition function, the Lower Confidence Bounds (Eq. 4) was used. In this equation, $\hat{\mu}_{\Theta}$ is the Gaussian Process estimation of the CV accuracy and $\hat{\Sigma}_{\Theta}$ is the covariance at parameter combination Θ

$$LCB(\Theta) = \hat{\mu}_{\Theta} - \hat{\Sigma}_{\Theta}. \quad (4)$$

The hyperparameter combination selected in the previous step was again evaluated using CV. Afterward, to refine the Gaussian Process and to determine the following combination, the respective tuple of hyperparameter and mean CV accuracy was added to the Gaussian Process. The procedure was repeated 25 times. The best hyperparameter combination was chosen to train the final model.

Model Training

During hyperparameter training and final model generation, XGBoost models, RFs, radial SVMs, polynomial SVMs, DTs, and LR models were trained. The preprocessing pipeline included centering, scaling, and median imputation. The entire preprocessing pipeline was implemented within the CV to avoid over-optimistic performance estimations [58]. The parameters were calculated for the CV training set and reused for the test and external datasets. The preprocessing was implemented using the Python package scikit-learn v0.23.2 [57].

Ensemble-based black-box XGBoost [4] models follow the idea of gradient boosting models [32]. It means that the

Table 5 Hyperparameters and intervals used to train the ML models

Model	Hyperparameter	Minimum	Maximum
LR	C	10^{-4}	10^2
	penalty	l2	none
DT	criterion	gini	entropy
	splitter	best	random
	max_depth	1	100
	min_samples_split	0	1
RF	n_estimators	250	1,250
	max_features	2	# features
	min_samples_leaf	1	20
XGBoost	n_estimator	1	500
	max_depth	1	20
	learning_rate	10^{-10}	1
	gamma	0	20
	min_child_weight	1	30
	subsample	0	1
	colsample_bytree	0	1
polynomial SVM	C	10^{-4}	10^2
	degree	1	10
	gamma	scale	auto
radial SVM	C	10^{-4}	10^2
	gamma	scale	auto

combination of multiple weak classifiers results in a strong, joint classifier. By learning the gradients of the previous classifier, gradient boosting fulfills this assumption. The final prediction consisted of the sum of weak classifier predictions. XGBoost is distributed as an open-source software library, and the main advantages are scalability, parallelization, and distributed execution. The hyperparameters and intervals used during Bayesian optimization are summarized in Table 5. The hyperparameter `n_estimators` sets the number of boosting iterations, `learning_rate` was the learning rate that preferences weak classifiers at early iterations, the minimum loss reduction required to split a node is defined by `gamma`, the hyperparameter `max_depth` sets the maximum depth of an individual tree, and the minimum number of observations in a child node was denoted as `min_child_weight`, `subsample` and `colsample_bytree` set the proportion of randomly subsampled training instances and features per iteration. The Python package `xgboost v1.2.0` [59] implemented the XGBoost algorithm.

RF [5] training was implemented using the Python package `scikit-learn v0.23.2` [57]. The RF algorithm is based on multiple DTs. Each DT was trained using randomly chosen features and subjects. Those subjects were selected using bootstrap sampling [60] on the training dataset. RF inference was computed by summarizing the individual DTs using a majority voting. The RF hyperparameters are summarized in Table 5. `n_estimators` sets the number of DTs, each

split used a random subset of `max_features` features, and the hyperparameter `min_samples_leaf` describes the minimum number of samples in a leaf node.

Support Vector Machines (SVMs) [9] were implemented using the Python package `scikit-learn v0.23.2` [57]. SVMs separate two classes using a decision boundary which was referred to as an n -dimensional hyperplane. Here, n is the number of features. To increase the robustness of the hyperplane for unknown observations, SVMs select the hyperplane with the largest distance from the observations. For this reason, the distance between the hyperplane and the observations was maximized using the hinge loss function [61]. The support vectors describe the observations closest to the hyperplane. Removing support vectors from the dataset directly influences the hyperplane. The cost parameter C enables SVMs to avoid overfitting, the higher C , the less complex an SVM is. Kernel functions help to model complex interactions. In this research, a polynomial and a radial kernel were implemented. The `degree` hyperparameter of the polynomial kernel controls the degree of the kernel, and high values lead to more complex hyperplanes. The `gamma` hyperparameter constraints the influence, and a single observation has on the hyperplane. If `gamma = scale`, $\frac{1}{\#features \cdot \sigma}$ was used as a value of `gamma`, if `gamma = auto`, a value of $\frac{1}{\#features}$ was used. The SVM hyperparameters and their ranges are summarized in Table 5.

In contrast to the black-box models, DTs [62] and LR models were selected as simple and interpretable comparison models. DTs were implemented using the Python package `scikit-learn v0.23.2` [57]. A DT consists of successively learned decision rules of the form $x \leq t$ for numerical or $x \in t$ for categorical features t is a threshold or a subset of values. The next decision rule was selected by the `splitter` which ranked all possible rules using a `criterion`. Decision rules were iteratively expanded until a maximum depth of `max_depth` was met or a percentage `min_samples_split` of samples were in a split.

LR [63] is a Generalized Linear Model (GLM) with a logistic link function. This link function allows the processing of binomial output variables. The logistic model function is given in Eq. 5. The model predicts the probability $P(Y = 1|X = x, \Theta)$ of observation x with given parameters Θ being in the positive class $Y = 1$. The LR algorithm was implemented using the Python package `scikit-learn v0.23.2` [57]

$$P(Y = 1|X = x, \Theta) = \frac{1}{1 + \exp(x \cdot \Theta)}. \quad (5)$$

Model Interpretation with Shapley Values

There are multiple methods to interpret ML models. An overview can be found in [10]. For example, DTs and LR

models are interpretable by design. However, black-box models often outperform those interpretable models but the interpretation of black-box models is complicated. In this research, model-agnostic Shapley values were used. Shapley values are local models, which explain the predictions of individual observations and thus enable high clinical benefit and high adaption to the black-box model.

Shapley values [64] are affiliated with coalition game theory and aim to decompose the prediction of a subject into the contributions of each feature. For this aim, Shapley values are based on the additive linear explanation model shown in Eq. 6. For a subject x , the model prediction $f(x)$ is decomposed into the feature contributions Φ_j , a simplified representation of the feature values x_j , and the average model prediction Φ_0 . A binned binary feature representation was used for tabular data, with N being the number of simplified features

$$f(x) = \Phi_0 + \sum_{j=1}^N \Phi_j x_j'. \quad (6)$$

The idea of using Shapley values to explain black-box ML models is to fairly decompose the contribution of each feature for the subject's prediction. Due to this fairness, the sum of all Shapley values is equal to the difference between the average model prediction and the probability prediction of a subject. Equation 7 shows that Shapley values are defined as the average, weighted contribution, a simplified feature has in all subsets. For the exact calculation of a Shapley value Φ_i for a given subject and feature i , it is required to determine the contribution of this feature for all subsets S of the entire feature set F . The investigation of each subset S requires the retraining and evaluation of the black-box model $f_S(S)$. With the help of the model performances trained with $(f_{S \cup i}(S \cup i))$ and without $(f_S(S))$ the feature at interest i , their differences were calculated. The weighted average difference across subsets builds the Shapley value. The weighting depends on the relative number of features $|S|$ in subset S . High weights were assigned to subsets with few and many features. In this way, the estimation of the main individual effects and the total effects are supported

$$\Phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} (f_{S \cup \{i\}}(S \cup \{i\}) - f_S(S)). \quad (7)$$

However, the number of subsets increases exponentially with the number of input features, leading to high computational effort for the exact calculation of Shapley values. Kernel SHapley Additive exPlanations (SHAP) [65] avoid time-consuming repeated training and evaluation by estimating Shapley values. This algorithm is based on Local Interpretable Model-agnostic Explanations (LIME) [66] and was implemented using the Python package `shap v0.38.1` [27]. A

new dataset containing variants of the observation at interest is created by permuting selected features. An additive linear model (Eqn. 8) with x' is a simplified representation of the black-box input features and $g(x')$ is the explanation model was fitted to the generated dataset

$$g(x') = \Phi_0 + \sum_{i=1}^M \Phi_i \cdot x'_i. \tag{8}$$

The weights Φ_i of the explanation model estimates the SHAP values for each subject and each feature. For tabular data, the simplified features are binned binary feature representations that represent if the original feature value or a permutation was used.

SHAP force plots [67] explain the model prediction of individual subjects using Shapley values. Features with positive Shapley values show strong positive effects on the prediction and small negative Shapley values represent small negative effects. SHAP force plots can be found in Fig. 15.

SHAP summary plots [67] summarize the explanations for the entire training dataset. Each point visualizes the feature value of a subject and the associated Shapley value. The color of a point depends on the subject’s feature value. On the vertical axis, the features are ordered by the mean absolute Shapley values. The plots were limited to the top ten features. SHAP summary plots can be found in Figs. 2, 5, 6, and 8.

There are some reasons, including out-of distribution sampling during Shapley value approximation and not taking into account feature correlation, why Shapley values should be used with caution for black-box model interpretability [68]. Therefore, it is important to compare Shapley value results with other ML explanation methods, or to reduce or consolidate correlated features [69]. In this work, forward selection was implemented to reduce the number of correlated features in the dataset, Shapley values were compared to classical feature importance measurements (“Classification Model”), and correlated features are consolidated to aspects.

Evaluation

The models were evaluated for the ADNI test set and the external AIBL and OASIS datasets. The performance was measured using accuracy (ACC) (Eq. 9), balanced accuracy (BACC) (Eq. 10), F1-Score (F1) (Eq. 11), and Matthews correlation coefficient (MCC) (Eq. 12). Table 6 visualizes the contingency table used for the calculation of those scores. Providing multiple scores for evaluation increased the comparability to other research. In comparison to accuracy, which focuses on correctly classified cases, the F1-Score focuses on incorrectly classified cases. The macro-averaging F1-Score was calculated to address both, the diseased

Table 6 Contingency table for the classification between patients and controls

Prediction	Diagnosis	
	Patient	Control
Patient	True positive (TP)	False positive (FP)
Control	False negative (FN)	True negative (TN)

and the healthy subject classification. Balanced accuracy is based on both, sensitivity and specificity and thus is suitable to evaluate imbalanced class problems. The MCC returns a value between 0 and 1, and is also suitable to handle imbalanced datasets

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \tag{9}$$

$$BACC = \frac{\frac{TP}{TP+FN} + \frac{TN}{TN+FP}}{2} \tag{10}$$

$$F1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)} \tag{11}$$

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}. \tag{12}$$

Additionally, the Area under the Receiver-Operating Curve (AUROC), which models the relationship between the True-Positive Rate (TPR—Eq. 13) and the False-Positive Rate (FPR—Eq. 14) for different classification thresholds was computed for all models. AUROC is suitable to investigate classification tasks with imbalanced datasets

$$TPR = \frac{TP}{TP + FN} \tag{13}$$

$$FPR = \frac{FP}{TN + FP}. \tag{14}$$

Results

In the following, the experimental results are presented. The MRI features selected using forward selection and the performances achieved for CN vs. AD, CN vs. MCI, MCI vs. AD, and sMCI vs. pMCI classification were given. SHAP summary plots compared the models trained using different feature sets, validation datasets, and classification models. The results of SHAP summary plots are compared to natural RF- and XGBoost-based feature importance scores and

Table 7 Features selected by forward selection using different ML methods as base classifiers for CN vs. AD classification

LR	DT	RF
sum_entorhinal	sum_Amygdala	sum_Amygdala
sum_Amygdala	sum_entorhinal	diff_parstriangularis
ratio_lingual	sum_Hippocampus	diff_superiorparietal
sum_middletemporal	ratio_supramarginal	sum_lateralorbitofrontal
diff_Lateral.Ventricle		sum_medialorbitofrontal
ratio_entorhinal		
XGBoost	SVM poly	SVM radial
sum_Amygdala	sum_entorhinal	sum_Amygdala
sum_middletemporal	sum_inferiorparietal	sum_entorhinal
sum_entorhinal	diff_Cortex	diff_Cortex
diff_lateralorbitofrontal	sum_Amygdala	sum_VentralDC
	ratio_paracentral	

Feature selection was exclusively used to reduce the number of MRI features

permutation importance scores. The influence of feature interactions for Shapley values is investigated and SHAP force plots explain individual model predictions of interesting subjects.

Feature Selection

The MRI features selected during forward selection for CN vs. AD classification and different ML methods used as base classifiers are summarized in Table 7. In this research, feature forward selection was used to reduce the number of MRI features and the influence of correlated features.

For the CN vs. AD detection task, the RF, and the polynomial SVM chose five features, the XGBoost, the DT, and the radial SVM chose four features and the LR chose six features. Overall, the six methods chose 16 different features. The most important feature for the RF, the XGBoost, the DT, and the radial SVM was the sum of the left and right amygdalae. For the polynomial SVM and the LR, the most important feature was the sum of the entorhinal cortices. Both features were previously associated with AD detection [70–73]. Previous research also shows that most of the selected features are associated with atrophy patterns in AD [74]. All methods also selected at least one difference or ratio of the left and right cortical or subcortical areas. Those features describe the asymmetry of both hemispheres. Brain asymmetry measurements were associated with AD [47–50] and were also successfully applied for ML models in this field [75].

Table 8 Features selected by forward selection using different ML methods as base classifiers for CN vs. MCI classification

LR	DT	RF
sum_middletemporal	sum_insula	sum_insula
ratio_isthmuscingulate	diff_insula	diff_isthmuscingulate
diff_paracentral	sum_fusiform	sum_inferiorparietal
diff_Cerebellum.White.Matter		sum_Cerebellum.White.Matter
XGBoost	SVM poly	SVM radial
ratio_inferiorparietal	sum_lingual	sum_temporalpole
sum_CerebralWhiteMatter	sum_Hippocampus	sum_inferiortemporal
ratio_VentralDC	ratio_rostralmid-dlefrontal	sum_caudalanteriorcingulate
diff_caudalanteriorcingulate	CSF	sum_Lateral.Ventricle
	sum_caudalanteriorcingulate	diff_precentral
	diff_isthmuscingulate	diff_Amygdala
	diff_Cerebellum.White.Matter	
	ratio_isthmuscingulate	

Feature selection was exclusively used to reduce the number of MRI features

The rankings of the forward selection for CN vs. MCI detection and different base classifiers are given in Table 8.

For CN vs. MCI detection, the RF, XGBoost, and LR base classifiers chose four features, the DT chose three features, the polynomial SVM chose eight features, and the radial SVM chose six features. Overall, the six ML methods chose 25 different features. Thus, in comparison to the CN vs. AD classification, the ML models show less agreement about the selected features. Consequently, the feature which was selected first in the forward selection process differed in five out of six methods. For the RF and the DT, the sum of the insular cortices was selected, the XGBoost classifier chose the ratio of the inferior parietal lobule, the polynomial SVM selected the sum of the lingual gyri, the SVM with the radial kernel chose the sum of the temporal pole volumes and the LR selected the sum of the left and right middle temporal gyri. Those features were previously associated with AD progression [70–74, 76–78]. Similar to the CN vs. AD classification, all models selected at least one feature describing the asymmetry of the cortical and subcortical brain regions.

The forward feature selection results of the six ML models for MCI vs. AD classification are summarized in Table 9. Four of the six models, namely RF, XGBoost, SVM poly, and LR selected five features. The DT chose six different features and the radial SVM selected two MRI features. Overall, the six methods selected 22 unique features.

Table 9 Features selected by forward selection using different ML methods as base classifiers for MCI vs. AD classification

LR	DT	RF
sum_entorhinal	sum_Hippocampus	sum_Hippocampus
sum_precuneus	sum_cuneus	sum_Amygdala
sum_VentralDC	sum_posteriorcingulate	diff_entorhinal
diff_frontalpole	ratio_Putamen	sum_isthmuscingulate
diff_rostralanteriorcingulate	sum_Cortex	ratio_lateralorbitofrontal
	ratio_parstriangularis	
XGBoost	SVM poly	SVM radial
diff_Lateral.Ventricle	sum_inferiortemporal	sum_inferiortemporal
diff_Cortex	Brain.Stem	ratio_frontalpole
sum_Cortex	sum_entorhinal	
sum_pericalcarine	sum_precuneus	
sum_precentral	ratio_precuneus	

Feature selection was exclusively used to reduce the number of MRI features

The most important features were the sum of the left and right hippocampi for the RF and DT model, the difference of the lateral ventricles for the XGBoost model, the sum of the inferior temporal gyri for both SVMs, and the sum of the entorhinal cortex volumes for the LR. Those features were previously associated with AD detection [70–73, 79, 80].

The results of the forward selection for the sMCI vs. pMCI classification task are summarized in Table 10. Five features were selected by the RF model, the XGBoost model chose six features, the DT selected only one feature, both SVMs chose four features, and the LR selected three features. Overall, the six methods picked 19 unique features. Three methods, namely the RF, the DT, and the SVM with the radial kernel selected the sum of the left and right amygdalae as the most important feature. The forward selection with the XGBoost base model first picked the sum of the hippocampi. The polynomial SVM selected the sum of the left and right precuneus and the LR chose the sum of the inferior temporal gyri. Those features were previously associated with AD detection [70–73, 80, 81].

Classification Tasks

In the following, the classification performances achieved for the four classification tasks are elaborated. The hyperparameters which reached the best accuracies during CV and which were thus used during training of the final models are summarized in Table 11.

Table 10 Features selected by forward selection using different ML methods as base classifiers for sMCI vs. pMCI classification

LR	DT	RF
sum_inferiortemporal	sum_Amygdala	sum_Amygdala
diff_middletemporal		sum_inferioparietal
sum_precentral		sum_entorhinal
		sum_lateraloccipital
		diff_superioparietal
XGBoost	SVM poly	SVM radial
sum_Hippocampus	sum_precuneus	sum_Amygdala
diff_lateralorbitofrontal	sum_inferiortemporal	diff_Inf.Lat.Vent
Brain.Stem	sum_rostralanteriorcingulate	sum_precuneus
sum_caudalmiddlefrontal	ratio_rostralanteriorcingulate	diff_middletemporal
diff_precentral		
sum_postcentral		

Feature selection was exclusively used to reduce the number of MRI features

CN vs. AD

The results achieved for CN vs. AD classification are summarized in Table 12. The no information rates were 60.36 % for the independent ADNI test set, 86.27 % for AIBL, and 78.05 % for OASIS. CN was the most frequent class in all datasets.

The best accuracy during CV of 99.68 % \pm 0.74 was achieved for the DT trained with feature selection and FS-3. This model also reached a perfect classification for the ADNI test set. All models trained for the CN vs. AD task reached accuracies higher than the no information rate for the ADNI dataset. The best AIBL accuracy of 95.94 % was achieved for the XGBoost model trained for FS-3 and with feature selection. This model also reached the best F1-Score (91.48 %) and the best MCC (0.830) for the AIBL dataset. The best balanced accuracies of 95.45 % for the AIBL dataset were reached for both SVMs trained with feature selection and FS-3. The LR model trained with feature selection for FS-3 reached the best AIBL AUROC of 99.55 %. Overall, two models achieved AIBL accuracies smaller than the no information rate of 86.27 %. Those models were the DTs trained with feature selection for FS-1 and FS-2.

The best OASIS accuracy of 90.58 % was achieved for the polynomial SVM which was trained with feature selection and FS-3. For OASIS, four models achieved accuracies worse than the no information rate of 78.05 %. Three of those models were trained on FS-1 and with feature selection, namely, the RF, the DT, and the SVM. The last model

Table 11 Hyperparameters tuned for CN vs. AD, CN vs. MCI, MCI vs. AD, and sMCI vs. pMCI classification. Hyperparameters: LR: C; penalty, DT: criterion; max depth; min samples split; splitter, RF: max features; min samples leaf; n estimators, XGBoost: colsample

bytree; gamma; learning rate; max depth; min child weight; n estimators; subsample, SVM poly: C; degree; gamma, SVM radial: C; gamma

Feature	Hyperparameters			
	CN vs. AD	CN vs. MCI	MCI vs. AD	sMCI vs. pMCI
FS-1				
LR	yes	{ 76.096; 12 }	{ 0.073; 12 }	{ 10.047; 12 }
LR	no	{ 0.0297; 12 }	{ 0.034; 12 }	{ 0.029; 12 }
DT	yes	{ 100; 0.143; best }	{ 100; 0.375; best }	{ 49; 0.994; best }
DT	no	{ 49; 0.994; best }	{ 31; 0.476; best }	{ 94; 0.823; random }
RF	yes	{ 5; 4; 1250 }	{ 4; 8; 955 }	{ 5; 1; 1250 }
RF	no	{ 77; 4; 1250 }	{ 95; 1; 1250 }	{ 28; 1; 250 }
XGBoost	yes	{ 0.814; 3.551; 0.025; 8; 1.0; 459; 0.765 }	{ 0.899; 0.660; 0.000; 13; 11.710; 488; 1.0 }	{ 1.0; 20.0; 1.0; 20; 1.0; 500; 1.0 }
XGBoost	no	{ 0.924; 3.795; 0.202; 12; 10.938; 299; 1.0 }	{ 0.671; 15.195; 0.000; 14; 7.070; 136; 0.967 }	{ 0.244; 4.526; 0.010; 13; 10.171; 500; 0.508 }
SVM poly	yes	{ 962.766; 1; scale }	{ 25.770; 1; auto }	{ 8.965; 1; auto }
SVM poly	no	{ 3.253; 1; auto }	{ 1.481; 1; auto }	{ 13.996; 3; auto }
SVM radial	yes	{ 0.717; scale }	{ 0.331; auto }	{ 1.483; scale }
SVM radial	no	{ 1.772; scale }	{ 1.685; scale }	{ 1.144; scale }
FS-2				
LR	yes	{ 0.095; 12 }	{ 24.121; none }	{ 0.083; 12 }
LR	no	{ 0.013; 12 }	{ 0.082; 12 }	{ 0.020; 12 }
DT	yes	{ 75; 0.354; best }	{ 47; 0.098; best }	{ 49; 0.994; best }
DT	no	{ 49; 0.994; best }	{ 100; 0.487; best }	{ 100; 0.366; best }
RF	yes	{ 2; 6; 1250 }	{ 3; 11; 250 }	{ 2; 1; 270 }
RF	no	{ 53; 3; 1250 }	{ 152; 1; 1250 }	{ 56; 1; 1250 }
XGBoost	yes	{ 0.885; 5.554; 0.012; 7; 3.119; 331; 0.296 }	{ 0.995; 6.791; 0.000; 9; 15.455; 477; 0.875 }	{ 1.0; 20.0; 1.0; 20; 1.0; 500; 1.0 }
XGBoost	no	{ 0.446; 1.499; 0.086; 10; 9.243; 361; 0.720 }	{ 0.897; 8.254; 0.120; 14; 4.872; 112; 0.936 }	{ 0.903; 11.976; 0.004; 5; 7.337; 376; 0.485 }
SVM poly	yes	{ 188.250; 1; auto }	{ 1000.0; 1; auto }	{ 1000.0; 1; scale }
SVM poly	no	{ 13.996; 3; auto }	{ 7.171; 1; auto }	{ 184.588; 3; auto }
SVM radial	yes	{ 2.526; scale }	{ 0.727; auto }	{ 1.343; auto }
SVM radial	no	{ 1.315; auto }	{ 1.367; auto }	{ 1.165; scale }
FS-3				
LR	yes	{ 26.861; 12 }	{ 0.586; 12 }	{ 0.375; 12 }
LR	no	{ 4.893; 12 }	{ 0.609; 12 }	{ 0.0209; 12 }
DT	yes	{ 100; 0.010; best }	{ 69; 0.079; best }	{ 87; 0.146; best }
DT	no	{ 100; 0.010; best }	{ 69; 0.079; best }	{ 69; 0.079; best }
RF	yes	{ 5; 1; 1236 }	{ 4; 5; 1250 }	{ 5; 5; 1126 }
RF	no	{ 36; 1; 250 }	{ 152; 20; 250 }	{ 41; 1; 1250 }
XGBoost	yes	{ 0.296; 19.279; 0.000; 14; 16.725; 445; 0.501 }	{ 1.0; 20.0; 1.0; 20; 1.0; 500; 0.571 }	{ 0.702; 3.777; 0.0013; 11; 5.021; 385; 0.230 }
XGBoost	no	{ 0.296; 19.279; 0.000; 14; 16.725; 445; 0.501 }	{ 1.0; 2.266; 0.000; 20; 1.0; 500; 0.685 }	{ 0.527; 13.650; 0.000; 19; 16.963; 254; 0.698 }
SVM poly	yes	{ 2.729; 1; auto }	{ 12.554; 1; auto }	{ 2.360; 1; auto }
SVM poly	no	{ 13.996; 3; auto }	{ 62.015; 1; scale }	{ 58.631; 3; auto }
SVM radial	yes	{ 3.342; auto }	{ 0.667; scale }	{ 0.464; auto }
SVM radial	no	{ 10.047; auto }	{ 10.047; auto }	{ 1.341; auto }

reaching an accuracy worse than the no information rate was the DT trained with FS-2 and feature selection.

CN vs. MCI

The results achieved for CN vs. MCI classification are summarized in Table 13. The no information rate for this task was 62.64 % for the ADNI test set, 82.44 % for AIBL, and 97.39 % for OASIS. MCI was the most frequent class in the ADNI dataset, whereas, for AIBL and OASIS, CN was.

The results achieved for CN vs. MCI classification were worse than those for the CN vs. AD task. The best accuracy during CV of 90.21 % ± 2.72 was achieved for the XGBoost model trained for FS-3 and without feature selection. The best accuracy for the ADNI test set was 91.58 % reached for two models. Both models, the radial SVM and the XGBoost model, were trained for FS-3 and with feature selection. The latter model also reached the best ADNI balanced accuracy and ADNI F1-Score. Overall, none of the models reached an ADNI accuracy worse than the no information rate of 62.51 %.

The results achieved for AIBL and OASIS were worse than the ADNI results. The best AIBL accuracy was 68.95 % achieved for two DTs trained with forward feature selection for FS-1 and FS-2. These models also reached the best AIBL balanced accuracies, AIBL F1-Scores, and AIBL MCCs.

The best OASIS accuracy of 55.05 % was reached for the DT trained without feature selection for FS-3. For the CN vs. MCI classification, all models achieved accuracies worse than the no information rates for OASIS and AIBL.

MCI vs. AD

The MCI vs. AD classification results are summarized in Table 14. The no information rate was 71.85 % for the ADNI test set, 57.23 % for AIBL, and 91.24 % for OASIS. The most frequent class was MCI for ADNI and AIBL as well as AD for OASIS.

The best CV accuracy of 89.39 % ± 2.99 was achieved for the RF trained without feature selection and FS-3. For the independent ADNI test set, the best accuracy was 88.66 %, reached by the RF and LR models trained with feature selection and FS-3. The first of those models also reached the best ADNI AUROC of 95.50 %, whereas the second model achieved the best ADNI F1-Score (85.50 %), and ADNI MCC (0.712). None of the models reached an ADNI accuracy worse than the no information rate. However, the DT trained without feature selection for FS-1 as well as the XGBoost and the DT both trained for FS-2 and with feature selection exactly achieved the no information rate of

subjects. The third most important feature was the summed volume of the entorhinal cortices. Consistently with AD atrophy, small volumes (colored in blue) were associated with AD progression. The same applied to the sum of the inferior parietal lobules and the amygdalae. Similar to the FS-1 and FS-2 models, the FS-3 model also learned young age (colored in blue) was associated with AD, although the mean age of the ADNI-CN group was younger than the mean age of the ADNI-AD group.

As can be seen in Table 13, for the CN vs. MCI classification, FS-3 outperformed FS-1 and FS-2 for the ADNI and AIBL performance scores. The best accuracies for OASIS were reached for FS-3, whereas FS-2 models outperform those models for F1-Score, balanced accuracy, AUROC, and MCC.

For the MCI vs. AD task, which is summarized in Table 14, the same applied to all ADNI and AIBL scores. For the OASIS dataset, the best accuracy and F1-Score were reached by FS-1 and the best balanced accuracy, AUROC, and MCC for FS-3.

The results for the sMCI vs. pMCI classification are shown in Table 15. Those results show that FS-3 outperformed FS-1 and FS-2 for all ADNI scores. For the AIBL dataset, the best accuracy, balanced accuracy, F1-Score, and MCC were achieved for FS-2, whereas the best AUROC was reached for FS-3.

To indicate whether the differences in ADNI test accuracies between the three feature sets are statistically significant, a Friedman test [93] (p -value < 0.05) was executed. For this investigation, the results of Tables 7, 8, 9, and 10 are summarized, resulting in 48 observations per feature set (six different models, two feature selection methods, and four tasks). The p -value of $2.2 \cdot 10^{-16}$ indicated statistically significant differences between the feature sets. To identify, which feature sets differed from each other, a pairwise Wilcoxon signed-rank test (p -value < 0.05) with Bonferroni adjustment was executed. A summary of the results is given in Table 16. The results FS-3 achieved significantly differed from FS-1 and FS-2. The FS-1 and FS-2 results showed no statistically significant differences.

Reproducibility

In this work, all models were trained using the ADNI dataset. Data from AIBL and OASIS subjects were used to test model reproducibility.

For all classification tasks, most models achieved worse results for AIBL and OASIS in comparison to the independent ADNI test set. The AIBL accuracies are plotted against the ADNI accuracies for all previously described models in Fig. 3. Overall, the AIBL accuracies were worse than those achieved for ADNI. The CN vs. AD classification models

achieved the best accuracies for ADNI and AIBL. The worst AIBL accuracies were achieved for CN vs. MCI classification, where all models reached AIBL accuracies worse than the no information rate. For the remaining classification tasks, most models reached AIBL accuracies better than the no information rate. For the sMCI vs. pMCI classification, no correlation between ADNI and AIBL accuracies was observed.

In Fig. 4, the OASIS accuracies of all previously described models are plotted against their ADNI accuracies. Similar to the AIBL results, the overall OASIS accuracies were worse than those achieved for ADNI. The best results for OASIS were achieved for CN vs. AD classification. Those models mainly reached accuracies better than the no information rate. The OASIS no information rates for the remaining classification tasks were larger than 90 % and all classification models trained for the ADNI dataset performed worse. However, the most frequent classes in OASIS and ADNI differed from each other for those classification tasks. For the OASIS dataset, the worst accuracy was achieved for MCI vs. AD classification. Reasons for the worse OASIS performances were, for example, differing MRI protocols and differences in the subject selection process.

To compare the model predictions for the three datasets, SHAP summary plots were visualized for the individual datasets in Fig. 5. Those plots show the Shapley values for the RF trained with feature selection and FS-3, which was trained for CN vs. AD classification. For all three datasets, the three most important features were the cognitive test scores, and bad scores were associated with disease progression. Those cognitive test scores were followed by the volumetric features, of the amygdalae, medial orbitofrontal cortices, and pars triangularis, as well as the AGE, the number of APOE ϵ 4 alleles, and the number of education years in slightly differing orders. For all volumetric features, biologically plausible associations [70, 82–84] were learned. The number of education years was not available in AIBL and OASIS, and those scores are therefore colored in grey.

SHAP summary plots for the RF trained with feature selection for CN vs. MCI classification based on FS-1 are shown in Fig. 6. The figure contains subplots for all three datasets. Overall, the Shapley values for this model were asymmetric. The positive Shapley values show stronger amplitudes than the negative ones. One explanation for this behavior might be that the MCI class was more frequent in the ADNI training dataset. For the ADNI and AIBL dataset, the most important feature was the sum of the inferior parietal lobules followed by the age and gender. The model learned that small brain volumes, young age, and male gender increased the risk to develop MCI. The volumetric observations correspond to previous research [70, 82–84]. The volume of the inferior parietal lobules

was the second most important feature for the OASIS dataset. Age was the most important feature for OASIS and the second most important feature for ADNI and AIBL. The model learned young age was associated with disease progression. It can be noted in Table 2 that the mean age of CN subjects is older than the mean age of MCI subjects in the ADNI dataset but not in the AIBL (Table 3) and OASIS (Table 4) datasets. The differences observed in the datasets might cause problems in model reproducibility. The feature representing the years of education was in the fifth position for the ADNI dataset. That information was not available in OASIS and AIBL and was thus colored in grey. Consistently, this feature was the least important one for both datasets. Overall, the ranking of the feature importance differed for all models.

Classification Model

In this research, six ML models were trained to compare their results to each other. A line plot of the accuracies achieved for the independent ADNI test set dependently on the classification task and the ML model is shown in Fig. 7. For the sMCI vs. pMCI classification, it can be seen that the performance variance is smaller for RF and XGBoost models in comparison to the remaining ML models. In addition, the polynomial SVMs achieved worse results for this classification task. Overall, the DT models were often outperformed by RF and XGBoost classifiers. The LR models outperformed the DTs in many cases, except for the CN vs. MCI classification. Overall, no ML model outperformed the remaining models.

To indicate whether the differences in ADNI test accuracies between the ML methods are statistically significant, a Friedman test (p -value < 0.05) was executed. For this investigation, the results of Tables 7, 8, 9, and 10 are summarized, resulting in 24 observations per feature set (three feature sets, two feature selection methods, and four tasks). The p -value of 0.006 indicated statistically significant differences between the ML models. A pairwise Wilcoxon signed-rank test (p -value < 0.05) with Bonferroni adjustment was executed, to identify, which model performances differed from each other. However, the results, summarized in Table 17, show that there are no statistically significant differences between ML models.

To visualize ML model differences, Fig. 8 shows SHAP summary plots for all six models. All models were trained using FS-3 with feature selection to distinguish between sMCI and pMCI subjects. The feature selection results in slightly different features within all models. Overall, the Shapley values had the largest deviance for the DT and the SVM with a polynomial kernel, followed by the LR model and the RF. The most important feature for all models except for the RF and the radial SVM was the LDELTOTAL cognitive test score. For this test score, all

models associated small feature values (colored in blue) with disease progression. For the radial SVM and the RF, LDELTOTAL was the second most important feature. The most important feature in the RF model was the sum of the left and right amygdalae. The model learned that large brain volumes decreased the patient's risk to develop AD. This observation is biologically plausible [70, 82–84]. The sum of the amygdala volumes was the third most important feature in the DT and the radial SVM. The number of ApoE ϵ 4 alleles was the most important feature for the radial SVM. The model learned that ApoE ϵ 4 is an AD risk factor, and the presence of ApoE ϵ 4 alleles is associated with AD progression. The number of ApoE ϵ 4 alleles is the second most important feature for the DT, and the LR, the third most important feature for the XGBoost model and the polynomial SVM, and the fourth most important feature for the RF. In this comparison, all models except for the DT and the polynomial SVM had at least one asymmetry feature within its top ten features. The decision tree only depended on three features, namely the LDELTOTAL cognitive test score, the number of ApoE ϵ 4 alleles, and the hippocampus volume. Most associations, the models learned, were biologically plausible. The radial RF showed two features with a biologically implausible association [70, 82–84]. The model learned that high volumes of the lateral occipital sulci, as well as a high number of education years, are associated with disease progression. Those features are ranked as the ninth and tenth important features in this model. Surprisingly, the association of the education feature was also learned for the SVMs and the LR. For the polynomial SVM, the summed volumes of the rostral anterior cingulate cortices show a biologically implausible [70, 82–84] association. Overall, biological plausibility should only be expected for high-performing models.

As a comparison, Fig. 9 visualizes the natural feature importance for the RF and XGBoost models, and the log odd's ratios for the LR model (ordered by the absolute log odd's ratio), and Fig. 10 shows the permutation importance of all models. The most important features of all natural feature importance plots and all permutation importance plots correspond to the SHAP summary plots.

The Kendall's tau rank correlation [94] between feature rankings for all SHAP models, natural XGBoost and RF feature importances, absolute log odd's ratios of the LR model, and permutation importance of all models is shown in Fig. 11. Due to the forward feature selection, the different models are trained on slightly different MRI features and the correlation was calculated for pairwise complete observations. However, the socio-demographic data, the number of ApoE ϵ 4 alleles, and the cognitive test scores were used to train all models. As the features within a specific model are identical, first, the SHAP values, the permutation importance,

Table 15 CV and test results for sMCI vs. pMCI classification

Model	Feature selection	CV		ADNI			AIBL					
		ACC ($\bar{x} \pm \sigma$) (in %)	ACC (in %)	BACC (in %)	AUROC (in %)	F1 (in %)	MCC	ACC (in %)	BACC (in %)	AUROC (in %)	F1 (in %)	MCC
FS-1												
LR	yes	62.21 ± 5.47	65.28	63.59	67.68	63.47	0.287	64.29	59.38	54.17	56.25	0.265
LR	no	67.69 ± 5.81	66.67	64.84	74.22	64.70	0.316	67.86	64.58	59.38	64.15	0.333
DT	yes	64.73 ± 5.33	68.06	66.56	69.79	66.61	0.346	53.57	51.04	47.92	50.48	0.022
DT	no	64.92 ± 5.84	68.06	67.19	70.12	67.29	0.348	53.57	51.04	45.83	50.48	0.022
RF	yes	68.30 ± 5.90	68.75	67.66	71.13	67.78	0.361	67.86	63.54	59.38	61.99	0.350
RF	no	68.64 ± 6.21	70.14	69.38	76.41	69.49	0.392	67.86	63.54	53.65	61.99	0.350
XGBoost	yes	65.50 ± 5.78	68.75	67.66	73.45	67.78	0.361	57.14	57.29	58.85	56.92	0.144
XGBoost	no	68.31 ± 5.91	70.83	69.84	76.07	70.00	0.405	57.14	54.17	54.17	53.33	0.091
SVM poly	yes	64.80 ± 5.57	59.03	57.19	65.51	56.76	0.152	57.14	51.04	51.04	42.86	0.040
SVM poly	no	64.33 ± 5.01	61.81	58.59	70.00	56.33	0.213	64.29	60.42	63.54	59.06	0.251
SVM radial	yes	66.30 ± 5.85	65.28	63.59	70.21	63.47	0.287	67.86	63.54	54.17	61.99	0.350
SVM radial	no	67.58 ± 5.35	68.75	67.34	75.43	67.43	0.360	78.57	75.00	66.67	75.44	0.603
FS-2												
LR	yes	65.60 ± 5.68	62.50	61.09	69.63	61.03	0.230	60.71	57.29	71.88	56.19	0.167
LR	no	68.24 ± 5.48	69.44	67.81	76.05	67.86	0.376	71.43	67.71	68.23	67.25	0.427
DT	yes	67.29 ± 5.43	65.97	63.59	72.14	62.97	0.304	67.86	65.62	55.47	65.71	0.331
DT	no	65.93 ± 5.69	61.81	61.56	68.02	61.49	0.230	64.29	65.62	66.41	64.29	0.312
RF	yes	70.14 ± 6.24	67.36	66.25	70.66	66.35	0.332	71.43	69.79	66.67	70.05	0.409
RF	no	69.01 ± 5.53	70.14	69.06	77.54	69.21	0.390	67.86	64.58	57.81	64.15	0.333
XGBoost	yes	66.70 ± 5.77	68.75	69.06	74.15	68.68	0.379	60.71	60.42	69.79	60.26	0.207
XGBoost	no	68.77 ± 5.60	70.14	69.06	76.31	69.21	0.390	64.29	60.42	59.38	59.06	0.251
SVM poly	yes	60.59 ± 5.53	57.64	56.25	58.09	56.10	0.129	64.29	61.46	71.35	61.11	0.251
SVM poly	no	64.83 ± 4.78	66.67	63.91	71.02	62.88	0.325	67.86	65.62	66.15	65.71	0.331
SVM radial	yes	68.89 ± 5.47	65.28	63.91	70.37	63.91	0.288	71.43	68.75	66.15	68.89	0.411
SVM radial	no	68.12 ± 5.85	69.44	68.12	76.23	68.24	0.375	75.00	72.92	72.40	73.33	0.486
FS-3												
LR	yes	69.79 ± 5.97	74.31	73.44	79.63	73.63	0.476	60.71	58.33	63.54	58.10	0.177
LR	no	69.55 ± 5.40	71.53	70.31	79.63	70.49	0.419	71.43	67.71	64.58	67.25	0.427
DT	yes	68.70 ± 5.85	74.31	73.28	79.72	73.51	0.476	53.57	53.12	44.53	53.03	0.062
DT	no	66.83 ± 5.47	68.06	66.41	75.28	66.40	0.346	53.57	51.04	43.75	50.48	0.022
RF	yes	70.75 ± 5.94	71.53	71.09	77.02	71.13	0.423	60.71	59.38	63.02	59.42	0.190
RF	no	70.64 ± 5.76	70.83	70.00	79.60	70.14	0.405	67.86	64.58	56.25	64.15	0.333
XGBoost	yes	69.11 ± 5.78	73.61	72.97	80.14	73.09	0.463	53.57	53.12	54.69	53.03	0.062
XGBoost	no	69.36 ± 5.58	73.61	72.50	78.48	72.72	0.462	64.29	61.46	52.60	61.11	0.251
SVM poly	yes	66.94 ± 5.32	64.58	62.50	72.83	62.08	0.271	60.71	58.33	58.33	58.10	0.177
SVM poly	no	66.35 ± 4.77	66.67	64.06	73.75	63.23	0.323	64.29	62.50	66.67	62.57	0.258
SVM radial	yes	70.09 ± 5.46	75.00	74.38	79.11	74.51	0.491	64.29	62.50	61.98	62.57	0.258
SVM radial	no	70.68 ± 5.19	68.06	66.72	79.00	66.80	0.346	82.14	80.21	73.96	80.95	0.640

All models were trained on ADNI and validated for an independent ADNI test set and external AIBL dataset. The best results in each section are highlighted in bold. No information rates: ADNI test set: 55.56 %, AIBL dataset: 57.14 %

and the feature importances are compared for each individual model. The SHAP values of the RF model and the permutation importance of the RF have a correlation coefficient of 0.82. The natural feature importance of the RF is only moderately correlated to the permutation importance of the RF (0.45) and weakly correlated to the RF SHAP values (0.35). The XGBoost SHAP values showed a very strong correlation of 0.82 to the natural XGBoost feature importance and a moderate correlation of 0.41 to the XGBoost permutation importance. The DT selected three features in all methods leading to a perfect correlation between the DT SHAP values, and the permutation importance as well as a very strong correlation of 0.89 for the DT SHAP values and the natural

feature importances. The SHAP values of the polynomial SVM showed a strong correlation to the permutation importance (0.67) of the same model. A moderate correlation of 0.53 was reached for the SHAP values of the radial SVM and the permutation importance of the same model. The SHAP LR values are strongly correlated (0.73) to the permutation importances and very strong correlated to the log odds (0.96).

As previously mentioned, the features within the different ML models differed which makes the comparison of inter- and intra-model correlations difficult. Considering the inter-model correlations, a perfect correlation of 1 was reached between the SHAP values of the RF and the SHAP values of

the XGBoost model, as well as the permutation importance of the LR and the SHAP values of the polynomial SVM.

The execution times of the different ML models, the SHAP algorithm, and the permutation importance calculation are summarized in Table 18. All experiments were executed on an NVIDIA® DGX-1⁵ supercomputer. The execution environment was an NVIDIA®-optimized⁶ Docker⁷ [95] container, running a Deepo⁸ image. The results showed that, except for the RF model, all models were trained in less than 1 s. The mean training time during CV was 5.66 s for the RF. The RF model was trained using twelve features. The long training time of the RF model was also reflected in the SHAP algorithm and the permutation importance. To compute the SHAP values of one subject, the RF model consumes 47.79 s, whereas the radial SVM which achieved the second slowest time requires only 13.51 s per subject. This results in an execution time of approximately 10 h to calculate the RF SHAP values of the entire dataset. The execution time for permutation importance was approximately 1 h for the RF model and 8 min for the XGBoost model reaching the second-longest execution time. Overall, it has to be mentioned that SHAP value calculation is a time-intensive process. However, the times presented can only be used as an orientation, and optimization is possible by for example clustering the background subjects of the SHAP algorithm. In this work, the samples of the entire training dataset were used as background subjects. The SHAP execution time depended on the number of features in a dataset, the time needed for model inference, the number of background subjects, and the number of subjects that should be explained.

Feature Dependency and Shapley Values

As feature correlations reduce the validity of explainability methods [96, 97], the previously explained SHAP summary plots are all generated using feature selection to avoid strong feature correlations in the dataset. Feature correlations can also make explainability more difficult [97] and may lead to biologically implausible explanations. The original dataset without feature selection contains many correlated features. To compare the explanations for such a dataset [69] developed a method to consolidate correlated features to aspects and compute permutation and SHAP importances for those aspects. First, correlated features of the entire training dataset are identified using Spearman rank correlation

coefficients. Hierarchical agglomerative clustering [97] was used to create a dendrogram. In this work, a threshold of $H = 0.5$ determining the least correlated features in a group filtered the resulting aspects from the dendrogram. The permutation and SHAP importances are computed by jointly permuting all features in an aspect. This work uses the python package dalex v1.4.1 [98] for implementation.

The resulting aspects computed for sMCI vs. pMCI classification and FS-3 without feature selection are shown in Table 19. The 161 features of FS-3 are consolidated to 79 aspects. Of those aspects, 14 included an individual feature. Of the remaining 65 aspects, nine included more than two features. As was expected, the differences and ratios of the same region are often correlated. At least one pair of ratio and difference for the same region was included within 49 aspects. Aspect_34 included four regions within the medial temporal lobe. Previous research showed, that those regions are important for the detection of AD progression [91, 92]. Aspect_30 consolidated three ventricular regions. Previous research found that ventricular enlargement was associated with AD progression [85, 86]. Aspect_46 included the cognitive test scores LIMMTOTAL and LDELTOTAL, and aspect_45 included the eTIV and the gender.

Using those aspects, the SHAP importances visualized in Fig. 12 were computed for the sMCI vs. pMCI classification without feature selection and for all ML models. The most important aspect for the RF, the XGBoost, and the DT was aspect_34, which consolidated the entorhinal cortices, the parahippocampal gyri, the amygdalae, and the hippocampi. Those brain areas were associated with AD in previous research [91, 92]. Aspect_34 also was the second most important aspect of the polynomial SVM and the third most important aspect in the LR and radial SVM. The most important aspect for the LR was aspect_27 which consolidated volumes of the fusiform, the inferior temporal, and the middle temporal gyri. This aspect also reached the second rank for the RF, radial SVM, and XGBoost models as well as the third place for the polynomial SVM. Previous research [91, 99] showed that those regions are affected in early AD stages. Aspect_46 consolidated the LDELTOTAL and LIMMTOTAL cognitive test scores, and was the most important aspect for both SVMs. This aspect also achieved the third rank for the XGBoost, RF and DT models, and the second rank for the LR. Overall, the most important aspects of the different models seem to be similar for the ML models.

Figure 13 shows the aspect permutation importance plots of the previously described models. The most important aspects chosen by SHAP importance and permutation importance matched for the DT, and the radial SVM. The most important aspect of the RF and XGBoost models was aspect_46 which included the cognitive test scores LDELTOTAL and LIMMTOTAL. This aspect reached the third rank using the SHAP method for both models. For the LR,

⁵ DGX-1: <https://www.nvidia.com/en-us/data-center/dgx-1/>, Accessed 2022-05-01.

⁶ NVIDIA®-Docker: <https://github.com/NVIDIA/nvidia-docker>, Accessed 2022-05-01.

⁷ Docker: <https://www.docker.com/>, Accessed 2022-05-01.

⁸ Deepo: <https://github.com/ufoym/deepo>, Accessed 2022-05-01.

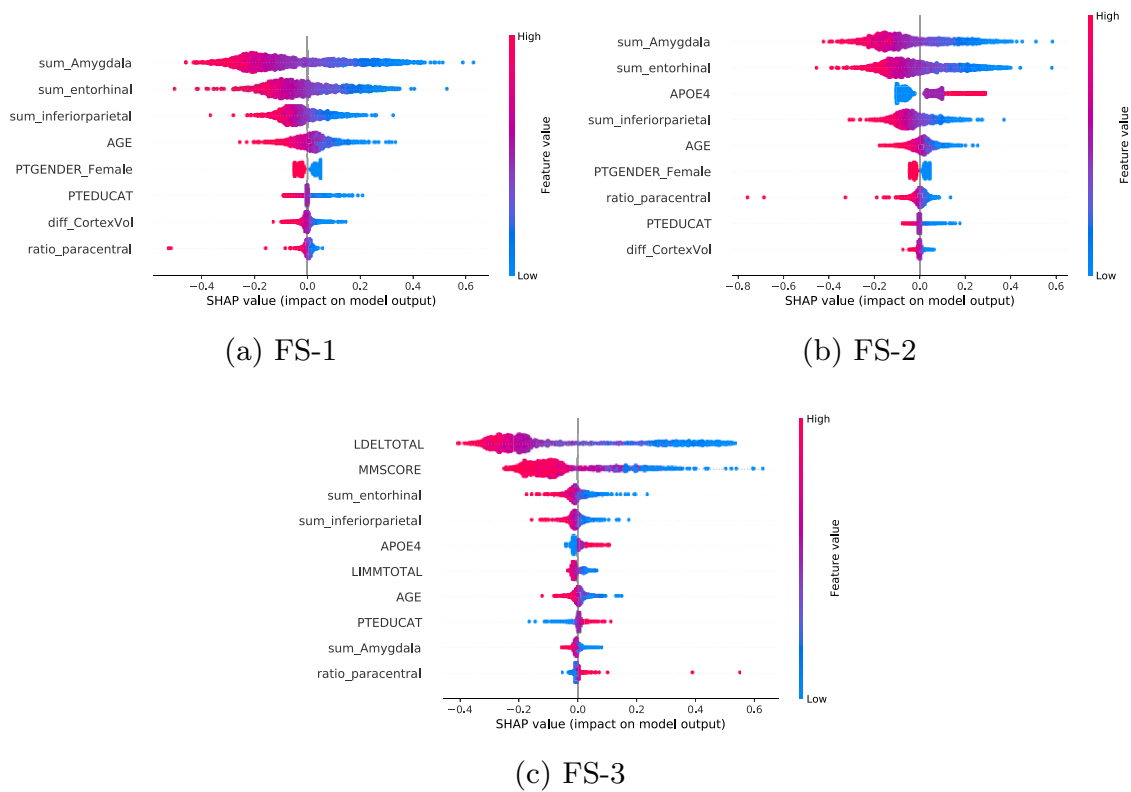


Fig. 2 SHAP summary plots of the polynomial SVM trained with feature selection for CN vs. AD. The plots visualize the Shapley values of $n=2,266$ subjects from the ADNI, AIBL, and OASIS datasets and the ten most important model features. Each subject is represented by a dot. The colors decode the subject’s feature expression.

High feature expressions are colored in red and small expressions in blue. The model learned that features with high Shapley values increased the patient’s AD risk. Each plot shows a model trained on a different feature set

Table 16 p -values of the pairwise Wilcoxon signed-rank test (p – value < 0.05) with Bonferroni adjustment to compare the differences in ADNI test accuracies between the three feature sets

	FS-1	FS-2	FS-3
FS-1	–	–	–
FS-2	0.95	–	–
FS-3	< 0.001	< 0.001	–

Statistically significant results are highlighted in bold

aspect_46 was also identified as the most important aspect. This aspect reached second place using the SHAP explanations. The highest permutation importance of the SVM was reached for aspect_34. This aspect reached second place using the SHAP method.

To investigate the correlation of the feature rankings, between the different methods, Kendall’s tau rank correlation between the SHAP feature ranking and the permutation method is visualized in Fig. 14. A very strong correlation of 1.00 was observed between the permutation importance and the SHAP values of the DT. A moderate correlation of 0.53 was observed between the SHAP values and the permutation

importance of the RF. The XGBoost SHAP rankings also showed a moderate correlation of 0.58 to the permutation importance of the same method. The SHAP values of the polynomial SVM are weakly correlated (0.36) to the model’s permutation importance. A very weak correlation of 0.07 was observed between the SHAP values of the radial SVM and their permutation importance rankings. The LR SHAP values are strongly correlated (0.74) to the permutation importance measurements of this model.

The inter-model correlations of the SHAP values showed a moderate correlation of 0.52 between the polynomial SVM and the radial SVM, as well as a strong correlation (0.65) between the SHAP values of the LR and the radial SVM SHAP values. The SHAP values of the RF was moderately correlated to the XGBoost SHAP values (0.52) and the SHAP values of the polynomial SVM. The SHAP values of the DT and LR model showed a very weak correlation (0.18). Overall, the SHAP values of the DT showed only weak correlations to the remaining ML models. The highest correlation of 0.39 was found for the SHAP values of the XGBoost model.

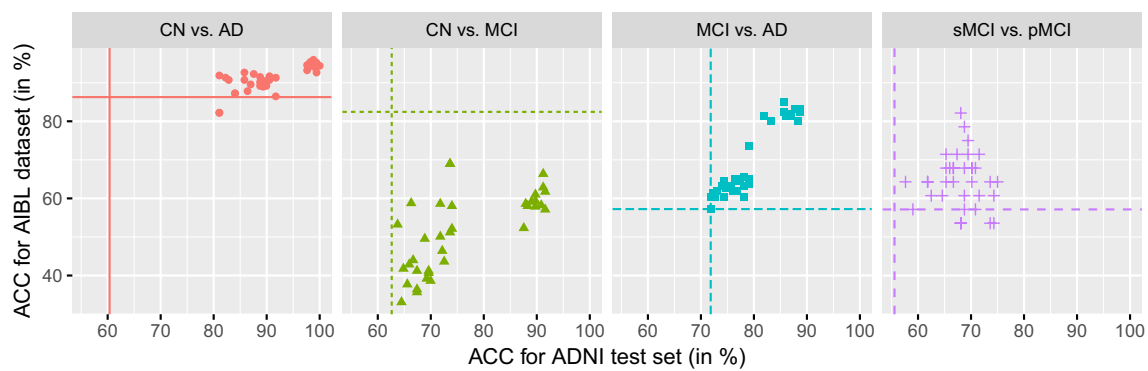


Fig. 3 Plot showing the accuracies achieved for the independent ADNI test set and the AIBL dataset for all models described in Tables 12, 13, 14, and 15. The no information rates for all classification tasks are visualized as horizontal lines for AIBL and as vertical lines for ADNI

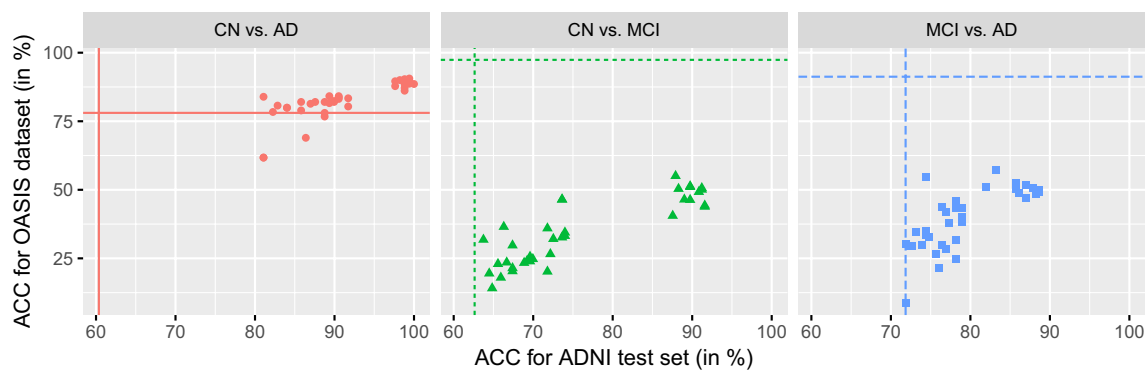


Fig. 4 Plot showing the accuracies achieved for the independent ADNI test set and the OASIS dataset for all models described in Tables 12, 13, 14 and 15. The no information rates for all classification tasks are visualized as horizontal lines for OASIS and as vertical lines for ADNI

Explanations of Individual Predictions

To investigate explanations for individual model predictions, Fig. 15 shows SHAP waterfall plots of four ADNI subjects. Those plots visualize the predictions for the RF trained with FS-3, and feature selection for sMCI vs. pMCI classification. SHAP waterfall plots explain the difference between the average model prediction value ($E[f(X)]$) and the subject's model prediction based on Shapley values. In all plots, the individual prediction was the probability of the subject being classified as pMCI. Features with a pathogenic expression are shown as red and protective expressions as blue arrows. The model prediction for the subject with PTID 027_S_1387 is explained in Fig. 15a. This is a subject from the ADNI test set and had a diagnosis of pMCI. The model prediction of this subject was 0.735. As this value was higher than 0.50, the subject was correctly classified as a pMCI subject. The most important feature with a pathogenic effect was the volume of the inferior parietal lobules. A relatively small normalized feature value of 0.237 was observed. The Shapley value of this feature was 0.12, and

thus, this feature expression increased the model prediction by 0.12. The LDELTOTAL cognitive test score reached a feature value of 3, which was a relatively bad test performance and thus increased the subject's risk to develop AD. Surprisingly, the relatively old age of 85.6 years decreased the patient's risk to develop AD by 0.03.

The model prediction for an sMCI subject (PTID: 037_S_4146) is demonstrated in Fig. 15b. This subject was sampled from the ADNI test set and reached a model prediction value of 0.149. The subject had a moderate-to-large volume of amygdalae which decreased the subject's risk of prospectively developing AD by -0.15 . The subject also has two ApoE ϵ 4 alleles, and as the presence of ApoE ϵ 4 alleles is a risk factor for AD, this increased the patient's risk. Additionally, the relatively high LDELTOTAL cognitive test score of 9 had a protective effect.

SHAP waterfall plots of subjects not included in the sMCI vs. pMCI dataset, because those pMCI subjects reverted to MCI at a later visit (explained in "Subject Selection") are visualized in Fig. 15c, d. The prediction of the subject with PTID 036_S_4430 is visualized in Fig. 15c. This MCI

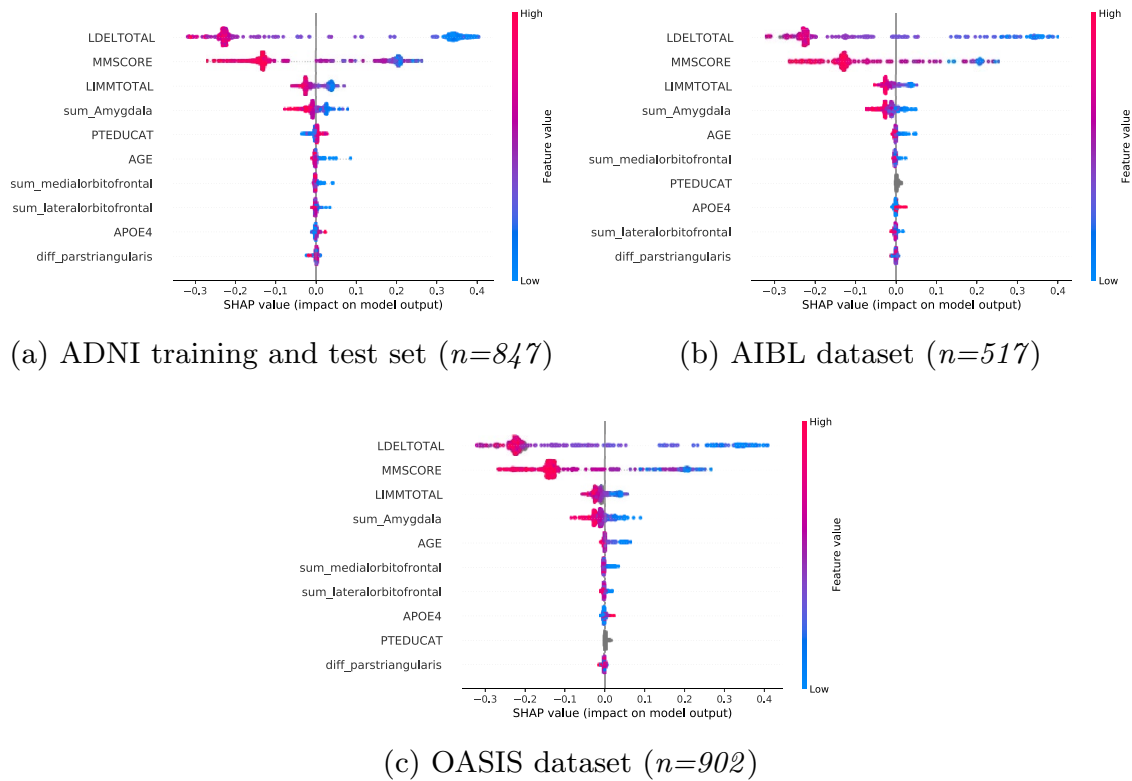


Fig. 5 SHAP summary plots of the RF trained with FS-3 and feature selection for CN vs. AD classification. The plots visualize the Shapley values of subjects from the ADNI, AIBL, and OASIS datasets and the ten most important model features. Each subject is represented by a dot. The colors decode the feature values of the subject. High

feature values are colored in red, whereas small feature values are colored in blue. The model learned that features with high Shapley values increased the patient’s risk to develop AD. Each plot shows the results on a different dataset

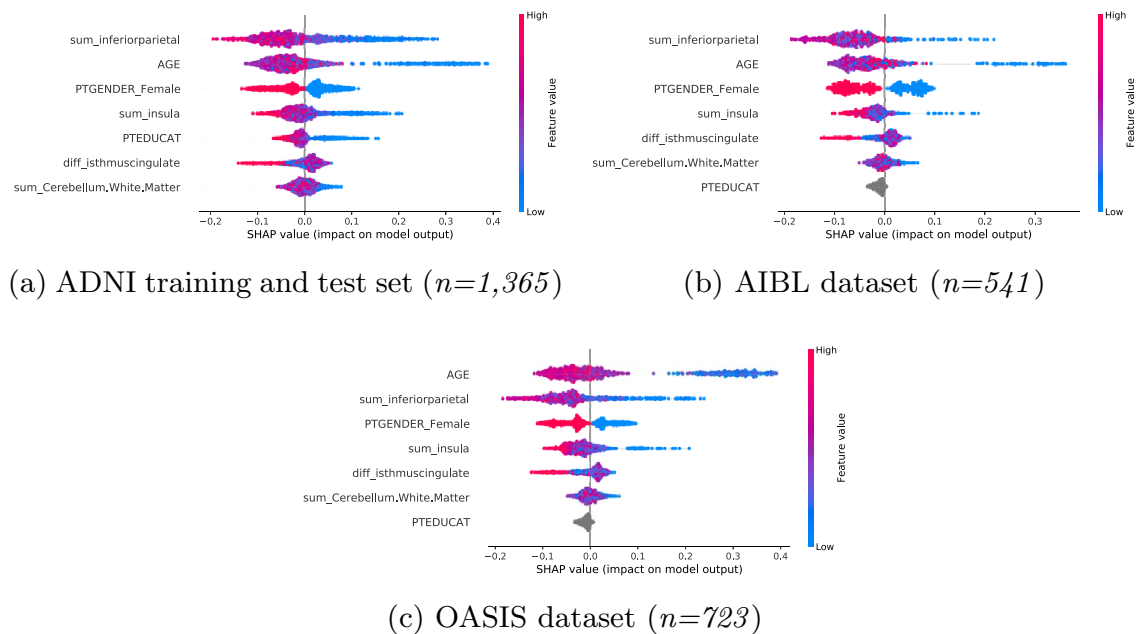


Fig. 6 SHAP summary plots of the RF trained with FS-1 and feature selection for CN vs. MCI classification. The plots visualize the Shapley values of subjects from the ADNI, AIBL, and OASIS datasets and the ten most important model features. Each subject is represented by a dot. The colors decode the feature values of the subject. High

feature values are colored in red whereas small feature values are colored in blue. The model learned that features with high Shapley values increased the patient’s risk to develop MCI. Each plot shows the results on a different dataset

Fig. 7 Line plot showing the accuracies achieved for the independent ADNI test set dependently on the classification tasks and the ML model. The plot includes all 216 models described in Tables 12, 13, 14 and 15. For each ML classifier, 36 models were included

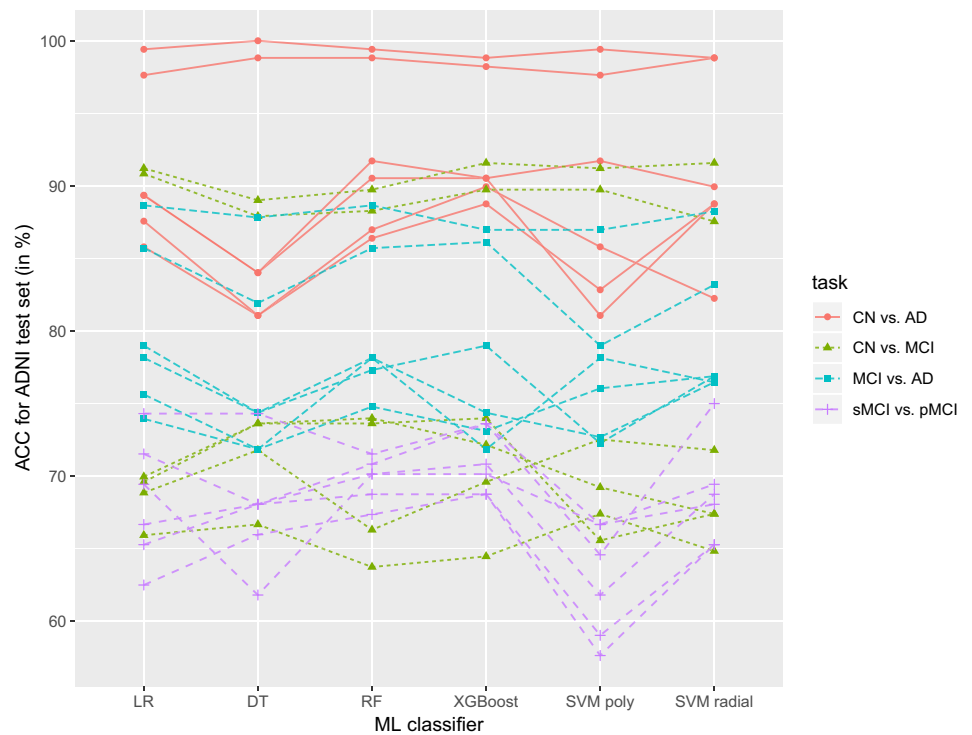


Table 17 *p*-values of the pairwise Wilcoxon signed-rank test (*p* – value < 0.05) with Bonferroni adjustment to compare the differences in ADNI test accuracies between the six ML models

	LR	DT	RF	XGBoost	SVM poly	SVM radial
LR	–	–	–	–	–	–
DT	0.622	–	–	–	–	–
RF	1.000	0.087	–	–	–	–
XGBoost	1.000	0.141	1.000	–	–	–
SVM poly	0.081	1.000	0.081	0.111	–	–
SVM radial	1.000	1.000	1.000	1.000	0.402	–

Statistically significant results are highlighted in bold

subject converted to AD 5.54 months after the baseline visit, but reverted to MCI 12.00 months after the baseline, and again converted to AD 23.64 months after the baseline visit. The last diagnosis for this subject was recorded after 83.54 months. The subject reached a model prediction of 0.707 and was thus classified as a pMCI subject. Additionally, the patient had a relatively small LIMMTOTAL cognitive test score of 2. The model learned that this poor test score increased the patient’s risk to develop AD by 0.09. Additionally, the subject had relatively small volumes of the amygdalae, which additionally decreased the patient’s risk to develop AD in the future. The AD risk of this patient was decreased by 0.03, because the subject has no ApoEε4 alleles.

The SHAP force plot for the subject with PTID 128_S_0135 is visualized in Fig. 15d. This MCI subject converted to AD 54.52 months after the baseline visit, reverted to MCI after 71.74 months and again converted to AD 83.90 months after the baseline visit which was also the last diagnosis available. However, in contrast to Fig. 15c, the subject reached a small model prediction value of 0.283 and was therefore classified as an sMCI subject. The most important factor decreasing the patient’s risk was the absence of ApoEε4 alleles. This factor decreased the model prediction by 0.06. Additionally, the LDELTOTAL cognitive test score of 8, which was relatively large, had a protective effect. The relatively small normalized volume of the lateral occipital sulci decreased the patient’s risk by 0.03. One reason for the classification

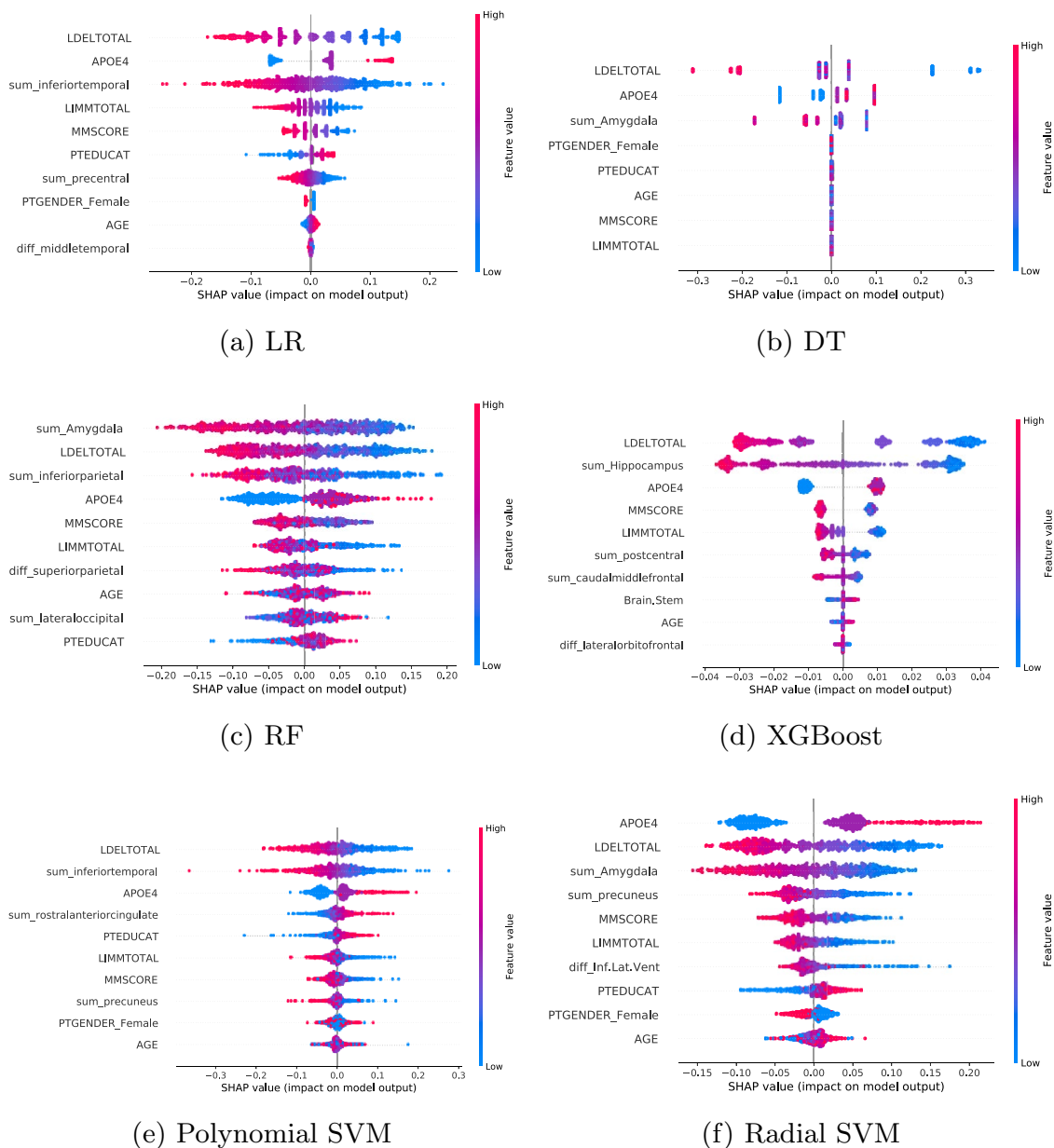


Fig. 8 SHAP summary plots for sMCI vs. pMCI classification trained using FS-3, feature selection, and multiple ML models. The Shapley values of $n=747$ ADNI and AIBL subjects and the ten most important model features are shown. Each dot represents a subject. The

colors decode the feature expressions. High expressions are colored in red and small ones in blue. The model learned that features with high Shapley values increased the patient’s AD risk

score might be that the conversion to AD was relatively late for this subject.

Discussion

In comparison to previous research [8], which exclusively trained tree-based models, this work trained several RFs, XGBoost models, DTs, SVMs, and LR models to detect different stages of AD. All models were trained using the

ADNI dataset and validated using independent test sets of the ADNI, AIBL, and OASIS cohorts. Bayesian optimization optimized for the best hyperparameters of the models. During this stage, CV was used to estimate the performance for independent test sets. The models were trained using three feature sets. The MRI features included summed volumes, differences, and ratios of predefined brain structures to investigate asymmetry structures associated with different AD stages. Forward feature selection was implemented to focus the models on the most important features and

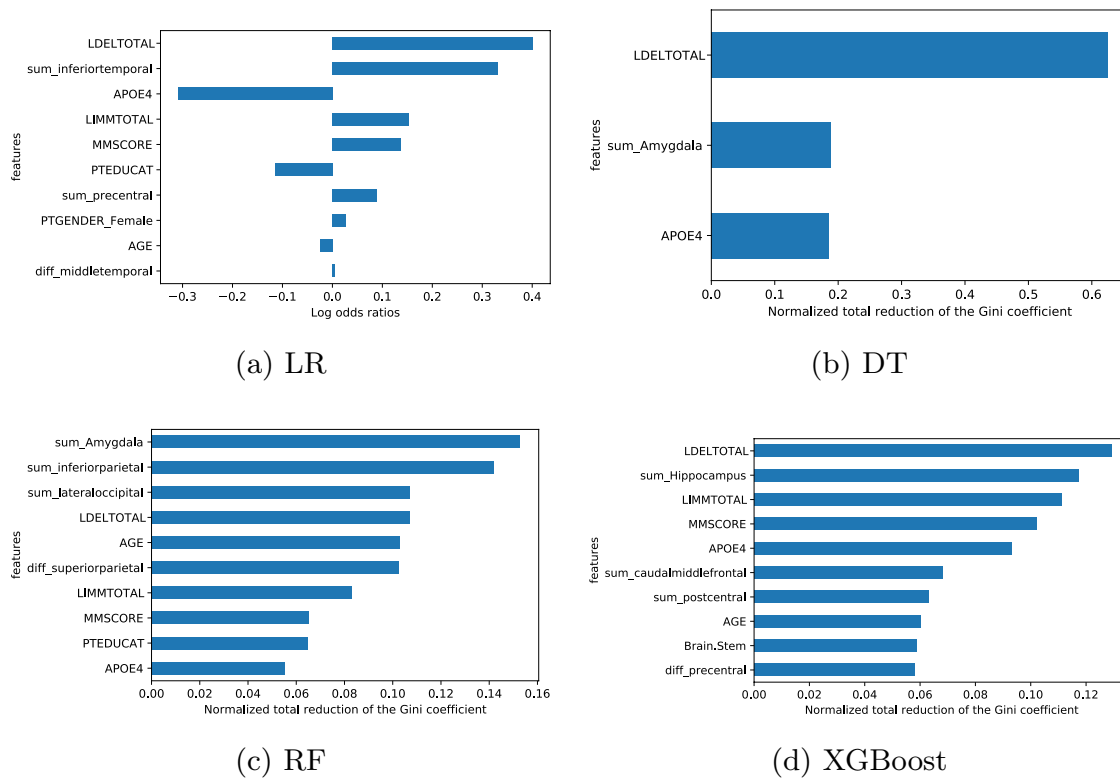


Fig. 9 Natural feature importance plots of the RF, XGBoost, and DT models and log odd’s ratios for LR models trained to distinguish between sMCI and pMCI subjects using FS-3 and feature selection. Each plot shows a different classification model

simultaneously avoid correlated features in the datasets. The performances of the different ML models as well as the different feature sets are compared to each other using Friedman tests and pairwise Wilcoxon signed-rank tests with Bonferroni adjustment. SHAP summary plots were used to visualize and interpret those models. The resulting Shapley values were compared to permutation importance of all models as well as natural feature importances of the RF and XGBoost models and to log odd’s ratios of the LR models. As correlated features reduce the validity of explainability methods like permutation importance and SHAP [96], those were also calculated consolidating correlated features to aspects. SHAP force plots investigated individual predictions of interesting subjects.

The experimental results showed that the forward feature selection chose brain regions that were previously associated with AD progression [70–73, 79, 80] for all classification tasks and models. The performances achieved for models trained with forward feature selection did not outperform the models trained on the entire feature set.

The pairwise Wilcoxon signed-rank tests with Bonferroni adjustment showed that the results of models trained with

FS-3, which included cognitive test results, outperformed those models trained for FS-1 and FS-2 for all classification tasks and the ADNI test set. The improvements for FS-3 models in comparison to FS-1 and FS-2 models were smaller for sMCI vs. pMCI than for the baseline classification tasks. The SHAP summary plots of all feature sets mainly showed biologically plausible associations and the most important features for the CN vs. AD classification using FS-3 and the polynomial SVM were the cognitive test scores LDELTOTAL and MMSCORE.

The results for the AIBL and OASIS test sets showed less clear advantages of FS-3. Reasons for this were, among others, differences in the subject recruitment process, leading to differences in socio-demographics and differing MRI protocols across studies. However, the CN vs. AD models were successfully transferred to AIBL and OASIS by mostly achieving classification accuracies better than the no information rate. Additionally, the models trained for MCI vs. AD classification and sMCI vs. pMCI classification were successfully transferred to the AIBL dataset. For CN vs. MCI classification, poor results worse than the no information rate were achieved for the AIBL and OASIS datasets.

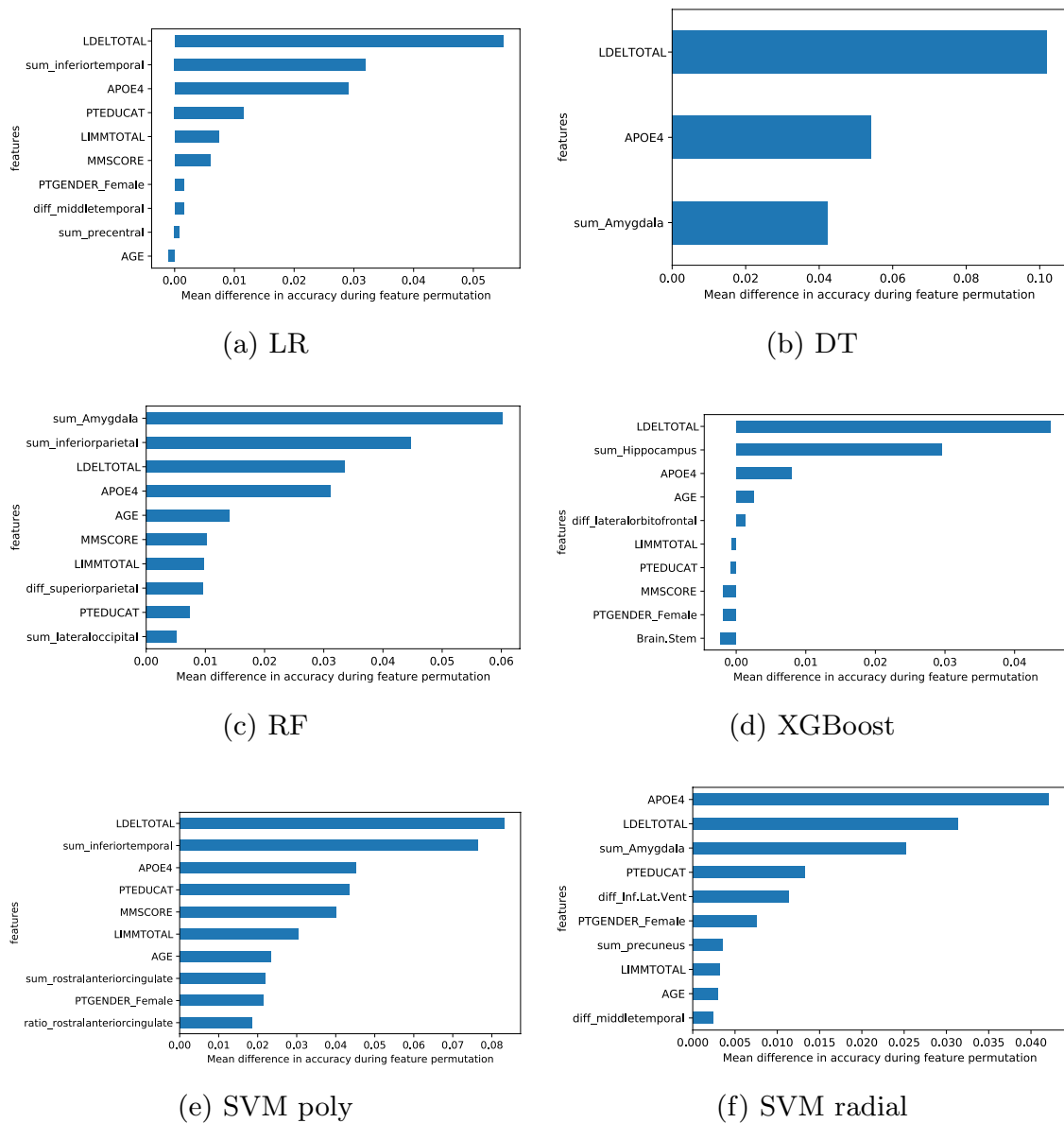


Fig. 10 Permutation importance plots of all six ML models trained to distinguish between sMCI and pMCI subjects using FS-3 and feature selection. Each plot shows a different classification model

However, the SHAP summary plots of those models mainly showed biologically plausible results. It was observed that age was a highly important feature in some of those models, which might cause problems transferring those models to datasets with differing demographic distributions. Poor results were also achieved for the MCI vs. AD classification and the OASIS dataset.

Some of the black-box models outperformed the simple and interpretable DTs. However, the pairwise Wilcoxon signed-rank tests with Bonferroni adjustment

($p - value < 0.05$) showed no significant differences. No model stood out among the black-box models. However, different ML models learned different associations which mostly were biologically plausible. The SHAP summary plots were compared to the permutation importance of all models, and natural RF and XGBoost feature importances, as well as absolute log odds of the LR and agreed for many features. The feature rankings of all models were compared to each other using Kendall’s rank correlation and showed moderate-to-strong correlation.

Fig. 11 Plot showing Kendall’s tau correlation between feature importances of all SHAP models, permutation importance of all models, and natural XGBoost and RF feature importances, as well as log odd’s ratios of the LR models for FS-3, feature selection, sMCI vs. pMCI classification and the ADNI and AIBL datasets ($n=747$)

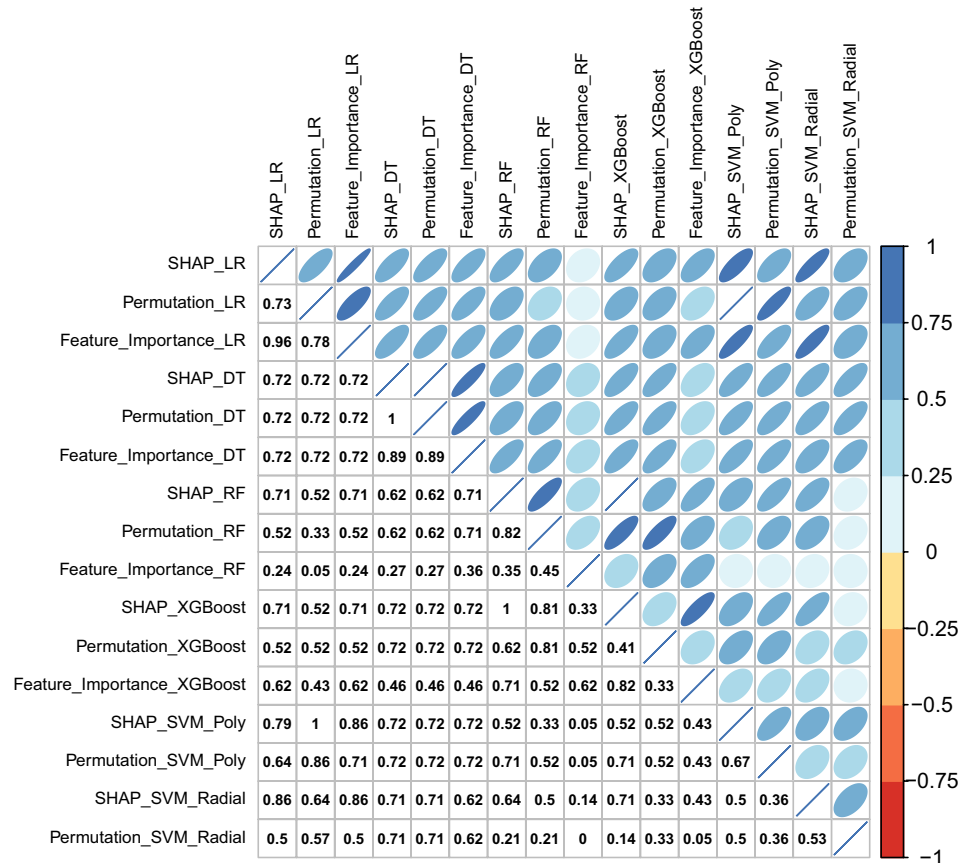


Table 18 Execution times of the different ML models and explainability methods for FS-3, feature selection, sMCI vs. pMCI classification, and the ADNI and AIBL datasets ($n=747$)

ML model	# model features	Mean CV training time (in h)	Mean SHAP time / subject (in h)	Global SHAP time ($n = 747$) (in h)	Global time permutation importance (in h)
LR	10	00:00:00.02	00:00:03.89	00:48:23.00	00:02:04.32
DT	8	00:00:00.01	00:00:01.18	00:14:41.00	00:01:39.61
RF	12	00:00:05.66	00:00:47.79	09:54:59.00	00:54:19.27
XGBoost	13	00:00:00.18	00:00:04.37	00:54:22.00	00:07:36.85
SVM Poly	11	00:00:00.15	00:00:08.09	01:40:44.00	00:05:05.77
SVM Radial	11	00:00:00.19	00:00:13.51	02:48:15.00	00:07:33.02

Because feature dependency structures reduce the validity of SHAP values and permutation importances, those feature importances were also computed by consolidating correlated features using aspects. The results show that the models depended on biologically plausible features. However, the feature rankings of the SHAP values and the permutation importances showed a weaker correlation than those calculated for the models with forward-selected features.

Individual predictions, which are important in clinical practice, were interpreted using SHAP waterfall plots.

Limitations

The approach proposed in this article had several limitations. First, both external datasets had a clear focus on CN subjects and were thus imbalanced which makes the interpretation of model generalizability hard. The external validation of the

Table 19 Aspects extracted with hierarchical agglomerative clustering for Spearman rank correlation and a threshold $H = 0.5$ for sMCI vs. pMCI classification and FS-3 without feature selection

Aspect	Features
aspect_1	[diff_Caudate, ratio_Caudate]
aspect_2	[diff_Putamen, ratio_Putamen]
aspect_3	[diff_inferiortemporal, ratio_inferiortemporal]
aspect_4	[diff parahippocampal, ratio parahippocampal]
aspect_5	[diff_entorhinal, ratio_entorhinal]
aspect_6	[diff_Lateral.Ventricle, ratio_Lateral.Ventricle]
aspect_7	[diff_Hippocampus, ratio_Hippocampus]
aspect_8	[diff_Inf.Lat.Vent, ratio_Inf.Lat.Vent]
aspect_9	[diff_temporalpole, ratio_temporalpole]
aspect_10	[diff_Amygdala, ratio_Amygdala]
aspect_11	[diff_posteriorcingulate, ratio_posteriorcingulate]
aspect_12	[diff_CerebralWhiteMatter, ratio_CerebralWhiteMatter]
aspect_13	[diff_Cortex, ratio_Cortex]
aspect_14	[diff_middletemporal, ratio_middletemporal]
aspect_15	[diff_lingual, ratio_lingual]
aspect_16	[diff_cuneus, ratio_cuneus]
aspect_17	[diff_pericalcarine, ratio_pericalcarine]
aspect_18	[diff_Thalamus.Proper, ratio_Thalamus.Proper]
aspect_19	[diff_Pallidum, ratio_Pallidum]
aspect_20	[diff_VentralDC, ratio_VentralDC]
aspect_21	[diff_Accumbens.area, ratio_Accumbens.area]
aspect_22	[sum_caudalanteriorcingulate, sum_rostralanteriorcingulate]
aspect_23	[sum_frontalpole, sum_lateralorbitofrontal, sum_medialorbitofrontal, sum_parsorbitalis, sum_rostralmiddlefrontal]
aspect_24	[sum_parsopercularis, sum_parstriangularis]
aspect_25	[sum_insula, sum_superiortemporal, sum_transversetemporal]
aspect_26	[sum_bankssts, sum_inferiorparietal]
aspect_27	[sum_fusiform, sum_inferiortemporal, sum_middletemporal]
aspect_28	[sum_precuneus, sum_superiorparietal, sum_supramarginal]
aspect_29	[sum_caudalmiddlefrontal, sum_paracentral, sum_postcentral, sum_precentral, sum_superiorfrontal, sum_Cortex]
sum_posteriorcingulate	[sum_posteriorcingulate]
aspect_30	[X3rd.Ventricle, sum_Inf.Lat.Vent, sum_Lateral.Ventricle]
sum_Accumbens.area	[sum_Accumbens.area]
AGE	[AGE]
CSF	[CSF]
sum_CerebralWhiteMatter	[sum_CerebralWhiteMatter]
aspect_31	[Brain.Stem, sum_Cerebellum.Cortex, sum_Cerebellum.White.Matter]
aspect_32	[sum_Thalamus.Proper, sum_VentralDC]
aspect_33	[sum_Pallidum, sum_Putamen]
aspect_34	[sum_entorhinal, sum parahippocampal, sum_Amygdala, sum_Hippocampus]
sum_temporalpole	[sum_temporalpole]
aspect_35	[sum_cuneus, sum_lingual, sum_pericalcarine]
sum_isthmuscingulate	[sum_isthmuscingulate]
sum_lateraloccipital	[sum_lateraloccipital]
aspect_36	[diff_fusiform, ratio_fusiform]
aspect_37	[diff_lateraloccipital, ratio_lateraloccipital]
aspect_38	[diff_supramarginal, ratio_supramarginal]
aspect_39	[diff_inferiorparietal, ratio_inferiorparietal]
aspect_40	[diff_superiorparietal, ratio_superiorparietal]
aspect_41	[diff_precuneus, ratio_precuneus]

Table 19 (continued)

Aspect	Features
aspect_42	[diff_bankssts, ratio_bankssts]
aspect_43	[diff_superiortemporal, ratio_superiortemporal]
aspect_44	[diff_parsorbitalis, ratio_parsorbitalis]
sum_vessel	[sum_vessel]
aspect_45	[EstimatedTotalIntraCranial, PTGENDER_Female]
sum_Caudate	[sum_Caudate]
aspect_46	[LIMMTOTAL, LDELTOTAL]
MMSCORE	[MMSCORE]
PTEDUCAT	[PTEDUCAT]
aspect_47	[diff_parstriangularis, ratio_parstriangularis]
aspect_48	[diff_parsopercularis, ratio_parsopercularis]
aspect_49	[diff_caudalmiddlefrontal, ratio_caudalmiddlefrontal]
aspect_50	[diff_rostralmiddlefrontal, ratio_rostralmiddlefrontal]
aspect_51	[diff_precentral, ratio_precentral]
X4th.Ventricle	[X4th.Ventricle]
APOE4	[APOE4]
aspect_52	[diff_postcentral, ratio_postcentral]
aspect_53	[diff_transversetemporal, ratio_transversetemporal]
aspect_54	[diff_rostralanteriorcingulate, ratio_rostralanteriorcingulate]
aspect_55	[diff_superiorfrontal, ratio_superiorfrontal]
aspect_56	[diff_caudalanteriorcingulate, ratio_caudalanteriorcingulate]
aspect_57	[diff_medialorbitofrontal, ratio_medialorbitofrontal]
aspect_58	[diff_frontalpole, ratio_frontalpole]
aspect_59	[diff_vessel, ratio_vessel]
aspect_60	[diff_lateralorbitofrontal, ratio_lateralorbitofrontal]
aspect_61	[diff_insula, ratio_insula]
aspect_62	[diff_isthmuscingulate, ratio_isthmuscingulate]
aspect_63	[diff_paracentral, ratio_paracentral]
aspect_64	[diff_Cerebellum.Cortex, ratio_Cerebellum.Cortex]
aspect_65	[diff_Cerebellum.White.Matter, ratio_Cerebellum.White.Matter]

sMCI vs. pMCI classification, which was medically more interesting than the baseline diagnoses, was based only on 28 AIBL subjects and no OASIS subjects. Future investigations should include more AD datasets knowing those cohorts differ in inclusion criteria. Possible cohorts might be the AD subset [14] of the HNR [15] or a subset of the National Alzheimer's Coordinating Center [100]. In this context, instead of diagnoses, different biomarkers should be addressed as endpoints. Another idea to increase the number of subjects in the datasets is to relax the exclusion criteria by also including subjects that reverted to MCI or CN, and use follow-up scans of subjects where the baseline scan failed for the MRI feature extraction pipeline. Due to the availability of data in the cohorts, and minimal invasive recording, only MRI, socio-demographics, the number of ApoE4 alleles, and cognitive test scores were included in

the investigations. However, PET scans and biomarkers have high medical relevance and should thus be considered in future investigations.

Although in comparison to previous research [8], the number of ML models was already increased, prospectively deep learning models like CNNs, which can automatically extract locally textural features from MRI scans should be investigated. However, currently, there is no consensus on whether those methods can improve AD detection. Much previous work in this area suffered from data leakage [23] or investigated the less challenging discrimination between AD and CN. The Bayesian optimization used for hyperparameter-tuning is a sequential method. Future work should therefore investigate the use of more effective parallelized methods such as presented in [101].

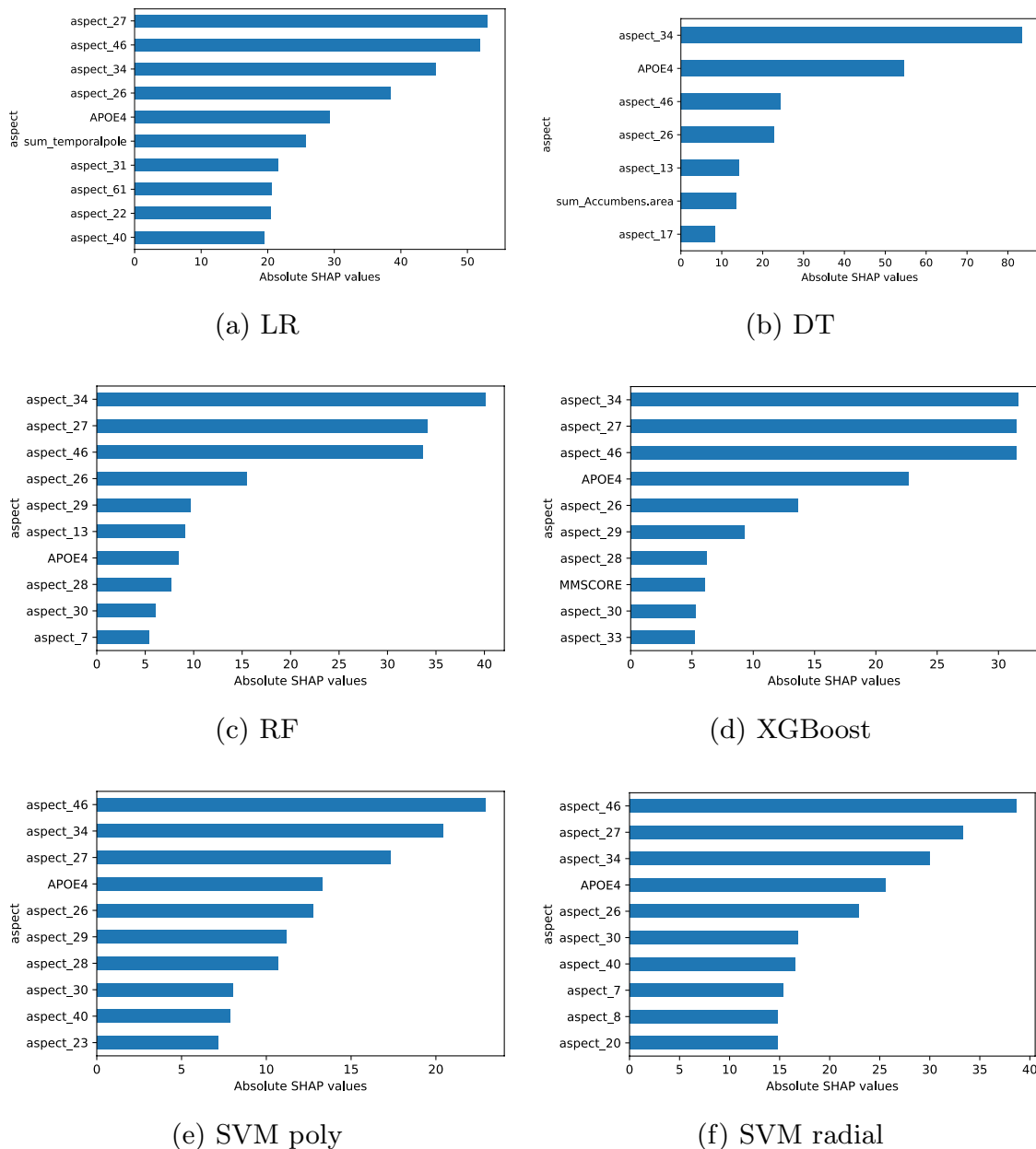


Fig. 12 SHAP aspect importance plots for all six ML models trained to distinguish between sMCI and pMCI subjects using FS-3 and no feature selection. Each plot shows a different classification model

Conclusion

This work extended a workflow [8] to explain ML black-box models trained to distinguish multiple AD stages using Shapley values. The differentiation of sMCI and pMCI subjects is of medical interest to recruit and monitor subjects for therapy studies. The approach was based on non-invasive features, including MRI volumes, socio-demographic data,

the number of ApoEε4 alleles, and cognitive test results. Volumetric features were extracted from the MRI scans using the FreeSurfer pipeline. The sum, difference, and ratio of the volumes of both hemispheres were calculated to investigate the brain asymmetry in multiple AD stages. Shapley sampling values were calculated to visualize the local feature associations of black-box RFs, XGBoost models, and SVMs. The experiments mainly showed biologically

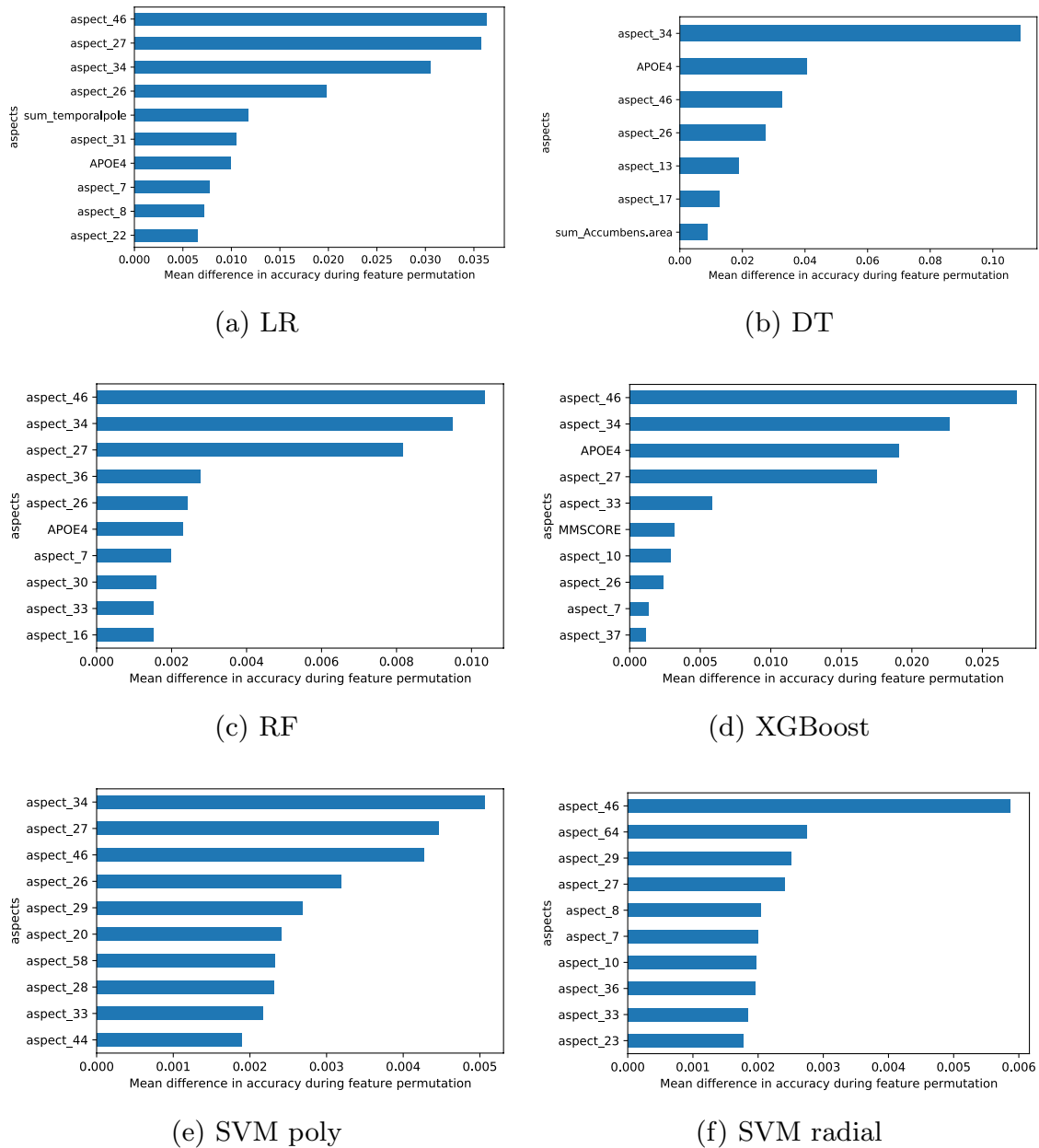


Fig. 13 Permutation aspect importance plots for all six ML models trained to distinguish between sMCI and pMCI subjects using FS-3 and no feature selection. Each plot shows a different classification model

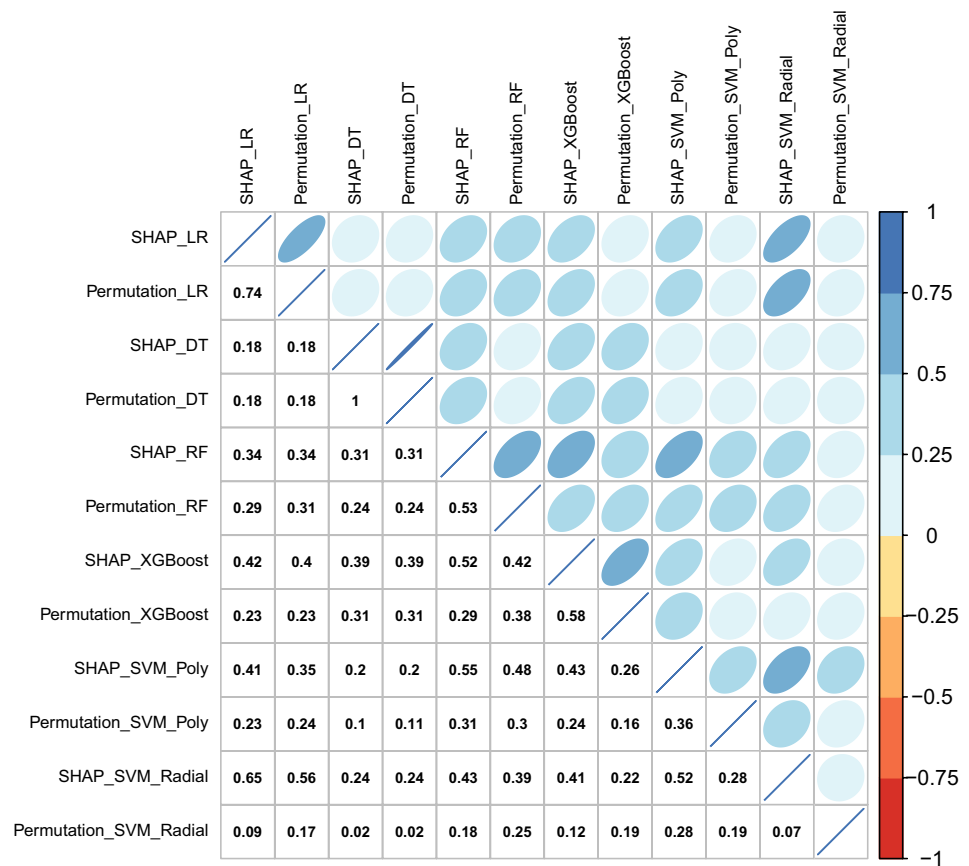
plausible associations and improved results for models including cognitive test scores. Those improvements were smaller for sMCI vs. pMCI classification.

For the investigation of model reproducibility, all models were trained for the ADNI dataset and validated for the external AIBL and OASIS cohorts. The ADNI models

achieved reasonable results for AIBL and CN vs. AD, MCI vs. AD, and sMCI vs. pMCI classification. For the OASIS test set, reasonable results were only reached for CN vs. AD classification.

Some of the performances of the black-box models outperformed the simple and interpretable DTs. None of the

Fig. 14 Plot showing Kendall's tau correlation between aspect importances of all SHAP models and permutation importance for FS-3, no feature selection, sMCI vs. pMCI classification, and the ADNI and AIBL data-sets ($n=747$)



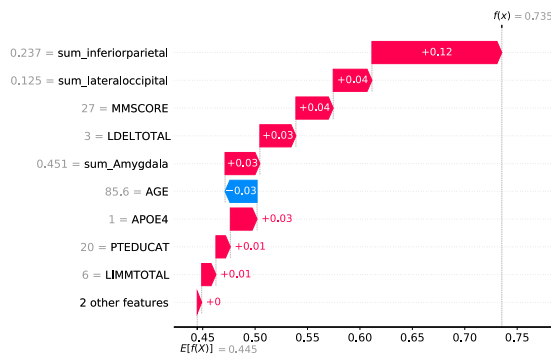
black-box models achieved outstanding results. SHAP summary plots were used to visualize the associations, the model learned between the features, and the AD diagnosis. The most important features of those plots were previously associated with AD progression. Additionally, those plots showed biologically plausible associations for most of the important features in all classification tasks.

SHAP force plots investigated individual model predictions. The comparison between SHAP values, natural and permutation feature importance showed moderate-to-strong correlations.

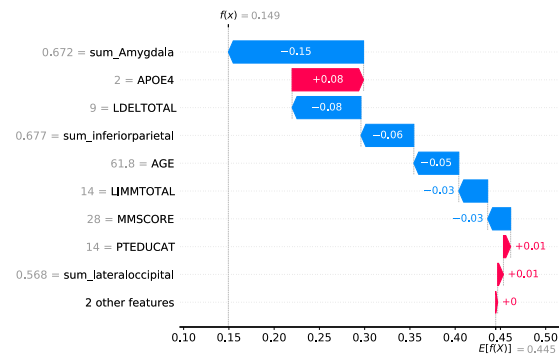
The investigation of the feature dependency structure and consolidating correlated features during the computation

feature importance computation for sMCI vs. pMCI classification showed that those models depended on features that were previously associated with AD.

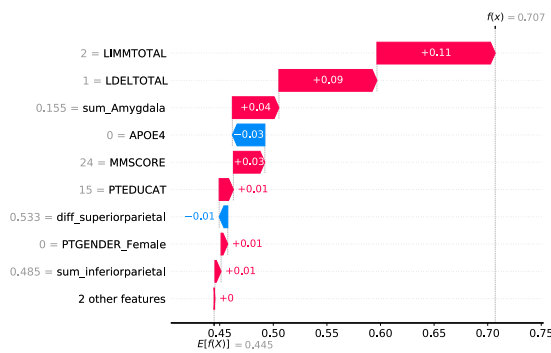
This work outperformed previous work [8] for the ADNI and AIBL classification results and CN vs. MCI classification and the AIBL results in MCI vs. AD classification for models trained without cognitive test scores. Additionally, the ADNI and AIBL results achieved for sMCI vs. pMCI classification trained with cognitive test scores outperformed the results of previous work.



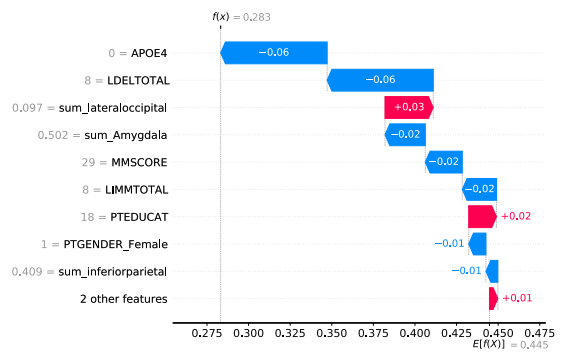
(a) PTID: 027_S_1387, diagnosis: pMCI, model prediction value: 0.735



(b) PTID: 037_S_4146, diagnosis: sMCI, model prediction value: 0.149



(c) PTID: 036_S_4430, diagnosis: pMCI subject which reverted to MCI, model prediction value: 0.707



(d) PTID: 128_S_0135, diagnosis: pMCI subject which reverted to MCI, model prediction value: 0.283

Fig. 15 SHAP waterfall plots for interesting ADNI subjects to explain the prediction of the RF trained with feature selection and FS-3 for sMCI vs. pMCI classification. Figure 15a, b shows subjects of the ADNI test set, and Fig. 15c, d explains model predictions of two subjects not included in the sMCI vs. pMCI dataset, because

those pMCI subjects reverted to MCI at a later visit. The arrow length indicates a Shapley value of the feature expression. Pathogenic feature expressions are shown as red and protective expressions as blue arrows. All volumetric features were scaled to a range between 0 and 1. Mean prediction value (base value): 0.445

Acknowledgements Data collection and sharing for this project was funded by the Alzheimer’s Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Ageing, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer’s Association; Alzheimer’s Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes

of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer’s Therapeutic Research Institute at the University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California. Data used in the preparation of this article were obtained from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) database (ADNI: <https://adni.loni.usc.edu>, Accessed: 2022-05-01). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data, but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found online (ADNI acknowledgement list: http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf, Accessed: 2022-05-01).

The authors thank Ahmad Idrissi-Yaghir, Department of Computer Science, University of Applied Sciences and Arts Dortmund,

44227 Dortmund, Germany, for the constructive proofreading of the manuscript.

Author Contributions The conceptualization of the study was carried out by CMF and LB. CMF and LB planned the experiments. LB implemented the software, executed the experiments, analyzed the data, and has written the original draft under the supervision of CMF. All authors read and approved the final manuscript.

Funding Open Access funding enabled and organized by Projekt DEAL. The work of Louise Bloch was partially funded by a PhD grant from University of Applied Sciences and Arts Dortmund, Dortmund, Germany.

Data Availability Statement Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (ADNI: <https://adni.loni.usc.edu>, Accessed: 2022-05-01) and Open Access Series of Imaging Studies (OASIS) (OASIS: <https://www.oasis-brains.org/>, Accessed: 2022-05-01). Details about data access are detailed there. The authors had no special access privileges others would not have to the data obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) or Open Access Series of Imaging Studies (OASIS) databases.

Code Availability Statement The workflow implementation will be available online after acceptance: <https://github.com/LouiseBloch/AlzheimerExplainableMLCorrelations>.

Declarations

Competing Interests The authors declare that they have no competing interests.

Ethics Approval Not applicable.

Consent to Participate The ADNI study was approved by the institutional review boards of the participating institutions. All participants gave informed written consent. More details can be found online (ADNI: <https://adni.loni.usc.edu>, Accessed: 2022-05-01). The AIBL study was approved by the institutional ethics committees of Austin Health, StVincent's Health, Hollywood Private Hospital and Edith Cowan University. All participants gave written informed consent before participating in the study. All OASIS participants consented to Knight Ageing and Disability Resource Center (ADRC)-related projects following procedures approved by the Institutional Review Board of Washington University School of Medicine. Participants consented to the use of their data by the scientific community and data sharing terms have been approved by the Washington University Human Research Protection Office.

Consent for Publication Consent for publication has been granted by ADNI administrators.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Alzheimer's Association. 2020 Alzheimer's disease facts and figures. *Alzheimers Dement.* 2020;16(3):391–460. <https://doi.org/10.1002/alz.12068>.
2. Patterson C. World Alzheimer Report 2018 - The State of the Art of Dementia Research: New Frontiers. Alzheimer's Disease International, London, Great Britain (2018). <https://www.alz.co.uk/research/WorldAlzheimerReport2018.pdf>, Accessed: 2021-10-10.
3. Sperling RA, Aisen PS, Beckett LA, Bennett DA, Craft S, Fagan AM, Ivatsubo T, Jack CR Jr, Kaye J, Montine TJ, Park DC, Reiman EM, Rowe CC, Siemers E, Stern Y, Yaffe K, Carrillo MC, Thies B, Morrison-Bogorad M, Wagster MV, Phelps CH. Toward defining the preclinical stages of Alzheimer's Disease: Recommendations from the national institute on aging - Alzheimer's association workgroups on diagnostic guidelines for Alzheimer's disease. *Alzheimers Dement.* 2011;7(3):280–92. <https://doi.org/10.1016/j.jalz.2011.03.003>.
4. Chen T, Guestrin C. XGBoost: A scalable tree boosting system. In: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), pp. 785–94. ACM, New York, United States 2016. <https://doi.org/10.1145/2939672.2939785>.
5. Breiman L. Random forests. *Mach Learn.* 2001;45(1):5–32. <https://doi.org/10.1023/A:1010933404324>.
6. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015;521(7553):436–44. <https://doi.org/10.1038/nature14539>.
7. Lapuschkin S, Wäldchen S, Binder A, Montavon G, Samek W, Müller K-R. Unmasking clever Hans predictors and assessing what machines really learn. *Nat Commun.* 2019. <https://doi.org/10.1038/s41467-019-08987-4>.
8. Bloch L, Friedrich CM. Developing a machine learning workflow to explain black-box models for Alzheimer's disease classification. In: Proceedings of the International Joint Conference on Biomedical Engineering Systems and Technologies (BIOSTEC 2021) - Volume 5: HEALTHINF, pp. 87–99. SciTePress, Setúbal, Portugal (2021). <https://doi.org/10.5220/0010211300870099>. INSTICC.
9. Cortes C, Vapnik V. Support-vector networks. *Mach Learn.* 1995;20(3):273–97. <https://doi.org/10.1007/BF00994018>.
10. Molnar C. Interpretable machine learning: A guide for making black box models explainable. <https://christophm.github.io/interpretable-ml-book/>, Accessed: 2021-10-10 2021.
11. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput.* 1997;9(8):1735–80. <https://doi.org/10.1162/neco.1997.9.8.1735>.
12. Rumelhart DE, Hinton GE, Williams RJ. Learning representations by back-propagating errors. *Nature.* 1986;323(6088):533–6. <https://doi.org/10.1038/323533a0>.
13. Pelka O, Friedrich CM, Nensa F, Mönninghoff C, Bloch L, Jöckel K-H, Schramm S, Sanchez Hoffmann S, Winkler A, Weimar C, Jöckel M. for the Alzheimer's Disease Neuroimaging Initiative: Sociodemographic data and APOE-ε4 augmentation for MRI-based detection of amnesic Mild Cognitive Impairment using deep learning systems. *PLoS ONE.* 2020;15(9):1–24. <https://doi.org/10.1371/journal.pone.0236868>.
14. Długaj M, Weimar C, Wege N, Verde PE, Gerwig M, Dragano N, Moebus S, Jöckel K-H, Erbel R, Siegrist J. Prevalence of mild cognitive impairment and its subtypes in the Heinz Nixdorf RECALL study cohort. *Dement Geriatr Cogn Disord.* 2010;30(4):362–73. <https://doi.org/10.1159/000320988>.
15. Schmermund A, Möhlenkamp S, Stang A, Grönemeyer D, Seibel R, Hirche H, Mann K, Siffert W, Lauterbach K, Siegrist J, Jöckel K-H, Erbel R. Assessment of clinically silent atherosclerotic

- disease and established and novel risk factors for predicting myocardial infarction and cardiac death in healthy middle-aged subjects: Rationale and design of the Heinz Nixdorf RECALL study. *Am Heart J.* 2002;144(2):212–8. <https://doi.org/10.1067/mhj.2002.123579>.
16. Petersen RC, Aisen PS, Beckett LA, Donohue MC, Gamst AC, Harvey DJ, Jack CR, Jagust WJ, Shaw LM, Toga AW, Trojanowski JQ, Weiner MW. Alzheimer's Disease Neuroimaging Initiative (ADNI). *Neurology.* 2010;74(3):201–9. <https://doi.org/10.1212/WNL.0b013e3181cb3e25>.
 17. Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017;618–26. <https://doi.org/10.1109/ICCV.2017.74>.
 18. Simonyan K, Vedaldi A, Zisserman A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *Computing Research Repository* 2014. [arxiv:1312.6034](https://arxiv.org/abs/1312.6034), Accessed: 2021-10-10.
 19. Springenberg JT, Dosovitskiy A, Brox T, Riedmiller M. Striving for simplicity: The all convolutional net. In: Proceedings of the International Conference on Learning Representations (ICLR) (workshop Track) (2015). <http://lmb.informatik.uni-freiburg.de/Publications/2015/DB15a>, Accessed: 2021-10-10.
 20. Zeiler MD, Fergus R. Visualizing and understanding convolutional networks. In: Fleet D, Pajdla T, Schiele B, Tuytelaars T (eds) Proceedings of the 13th European Conference on Computer Vision (ECCV), pp. 818–33. Springer, Basel, Switzerland 2014. https://doi.org/10.1007/978-3-319-10590-1_53.
 21. Yang C, Rangarajan A, Ranka S. Visual explanations from deep 3D convolutional neural networks for Alzheimer's disease classification. *Computing Research Repository* 2018. [arxiv:1803.02544](https://arxiv.org/abs/1803.02544), Accessed: 2021-10-10.
 22. Rieke J, Eitel F, Weygandt M, Haynes J-D, Ritter K. Visualizing convolutional networks for MRI-based diagnosis of Alzheimer's disease. In: Stoyanov D, Taylor Z, Kia SM, Oguz I, Reyes M, Martel A, Maier-Hein L, Marquand AF, Duchesnay E, Löfstedt T, Landman B, Cardoso MJ, Silva CA, Pereira S, Meier R (eds) Understanding and Interpreting Machine Learning in Medical Image Computing Applications, pp. 24–31. Springer, Basel, Switzerland 2018. https://doi.org/10.1007/978-3-030-02628-8_3.
 23. Wen J, Thibeau-Sutre E, Diaz-Melo M, Samper-González J, Routier A, Bottani S, Dormont D, Durrleman S, Burgos N, Colliot O. Convolutional neural networks for classification of Alzheimer's disease: Overview and reproducible evaluation. *Med Image Anal.* 2020;63: 101694. <https://doi.org/10.1016/j.media.2020.101694>.
 24. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y. Generative adversarial nets. In: Advances in Neural Information Processing Systems (NIPS), 2014; 2672–80. <https://papers.nips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afcc3-Paper.pdf> Accessed 2021-10-10.
 25. Wang X, Shen D, Huang H. Interpretable deep temporal structure learning model for early detection of Alzheimer's Disease. *bioRxiv* (2019). <https://doi.org/10.1101/2019.12.12.874784>
 26. Das D, Ito J, Kadowaki T, Tsuda K. An interpretable machine learning model for diagnosis of Alzheimer's disease. *PeerJ.* 2019;7:6543. <https://doi.org/10.7717/peerj.6543>.
 27. Lundberg SM, Lee S-I. A unified approach to interpreting model predictions. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R (eds) Advances in Neural Information Processing Systems vol. 30, pp. 4765–74. Curran Associates, Inc., New York, New York, US 2017. <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>, Accessed: 2021-10-10.
 28. Bloch L, Friedrich CM. Data analysis with Shapley values for automatic subject selection in Alzheimer's disease data sets using interpretable machine learning. *Alzheimers Res Ther.* 2021;13(1):155. <https://doi.org/10.1186/s13195-021-00879-4>.
 29. Ghorbani A, Zou J. Data Shapley: Equitable valuation of data for machine learning. In: Proceedings of the International Conference on Machine Learning (ICML), 2019;97:2242–51. <http://proceedings.mlr.press/v97/ghorbani19c/ghorbani19c.pdf> Accessed 2021-10-10.
 30. Cook RD. Detection of influential observation in linear regression. *Technometrics.* 1977;19(1):15–8. <https://doi.org/10.2307/1268249>.
 31. Hammond TC, Xing X, Wang C, Ma D, Nho K, Crane PK, Elahi F, Ziegler DA, Liang G, Cheng Q, Yanckello LM, Jacobs N, Lin A-L. β -Amyloid and tau drive early Alzheimer's disease decline while glucose hypometabolism drives late decline. *Commun Biol.* 2020. <https://doi.org/10.1038/s42003-020-1079-x>.
 32. Friedman JH. Greedy function approximation: A gradient boosting machine. *Ann Stat.* 2001;29(5):1189–232. <https://doi.org/10.1214/aos/1013203451>.
 33. Danso SO, Zeng Z, Muniz-Terrera G, Ritchie CW. Developing an explainable machine learning-based personalised dementia risk prediction model: A transfer learning approach with ensemble learning algorithms. *Front Big Data.* 2021;4: 613047. <https://doi.org/10.3389/fdata.2021.613047>.
 34. Börsch-Supan A, Brandt M, Hunkler C, Kneip T, Korbmayer J, Malter F, Schaaf B, Stuck S, Zuber o.b.o.t.S.C.C.T. Sabrina. Data resource profile: The Survey of Health, Ageing and Retirement in Europe (SHARE). *Int J Epidemiol.* 2013;42(4):992–1001. <https://doi.org/10.1093/ije/dyt088>.
 35. Ritchie CW, Ritchie K. The PREVENT study: A prospective cohort study to identify mid-life biomarkers of late-onset Alzheimer's disease. *BMJ Open.* 2012. <https://doi.org/10.1136/bmjopen-2012-001893>.
 36. Livingston G, Huntley J, Sommerlad A, Ames D, Ballard C, Banerjee S, Brayne C, Burns A, Cohen-Mansfield J, Cooper C, Costafreda SG, Dias A, Fox N, Gitlin LN, Howard R, Kales HC, Kivimäki M, Larson EB, Ogunniyi A, Orgeta V, Ritchie K, Rockwood K, Sampson EL, Samus Q, Schneider LS, Selbæk G, Teri L, Mukadam N. Dementia prevention, intervention, and care: 2020 report of the lancet commission. *The Lancet.* 2020;396(10248):413–46. [https://doi.org/10.1016/s0140-6736\(20\)30367-6](https://doi.org/10.1016/s0140-6736(20)30367-6).
 37. El-Sappagh S, Alonso JM, Islam SMR, Sultan AM, Kwak KS. A multilayer multimodal detection and prediction model based on explainable artificial intelligence for Alzheimer's disease. *Sci Rep.* 2021;11(1):2660. <https://doi.org/10.1038/s41598-021-82098-3>.
 38. Van Rossum G, Drake FL. Python 3 Reference Manual. CreateSpace, Scotts Valley, CA 2009. <https://www.python.org/>
 39. Ellis KA, Bush AI, Darby D, De Fazio D, Foster J, Hudson P, Lautenschlager NT, Lenzo N, Martins RN, Maruff P, Masters C, Milner A, Pike K, Rowe C, Savage G, Szoëke C, Taddei K, Villemagne V, Woodward M, Ames D. AIBL Research Group: The Australian Imaging, Biomarkers and Lifestyle (AIBL) study of aging: Methodology and baseline characteristics of 1112 individuals recruited for a longitudinal study of Alzheimer's disease. *Int Psychogeriatr.* 2009;21(4):672–87. <https://doi.org/10.1017/S1041610209009405>.
 40. LaMontagne PJ, Benzinger TL, Morris JC, Keefe S, Hornbeck R, Xiong C, Grant E, Hassenstab J, Moulder K, Vlassenko AG, Raichle ME, Cruchaga C, Marcus D (2019). OASIS-3: Longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and Alzheimer disease. *medRxiv.* <https://doi.org/10.1101/2019.12.13.19014902>

41. McKhann G, Drachman D, Folstein M, Katzman R, Price D, Stadlan EM. Clinical diagnosis of Alzheimer's disease. *Neurology*. 1984;34(7):939. <https://doi.org/10.1212/WNL.34.7.939>.
42. ...Winblad B, Palmer K, Kivipelto M, Jelic V, Fratiglioni L, Wahlund L-O, Nordberg A, Bäckman L, Albert M, Almkvist O, Arai H, Basun H, Blennow K, De Leon M, DeCarli C, Erkinjuntti T, Giacobini E, Graff C, Hardy J, Jack C, Jorm A, Ritchie K, Van Duijn C, Visser P, Petersen RC. Mild cognitive impairment - beyond controversies, towards a consensus: report of the international working group on mild cognitive impairment. *J Intern Med*. 2004;256(3):240–6. <https://doi.org/10.1111/j.1365-2796.2004.01380.x>.
43. Fischl B. FreeSurfer. *Neuroimage*. 2012;62(2):774–81. <https://doi.org/10.1016/j.neuroimage.2012.01.021>.
44. Desikan RS, Ségonne F, Fischl B, Quinn BT, Dickerson BC, Blacker D, Buckner RL, Dale AM, Maguire RP, Hyman BT, Albert MS, Killiany RJ. An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage*. 2006;31(3):968–80. <https://doi.org/10.1016/j.neuroimage.2006.01.021>.
45. Fischl B, Salat DH, Busa E, Albert M, Dieterich M, Haselgrove C, van der Kouwe A, Killiany R, Kennedy D, Klaveness S, Montillo A, Makris N, Rosen B, Dale AM. Whole brain segmentation: Automated labeling of neuroanatomical structures in the human brain. *Neuron*. 2002;33(3):341–55. [https://doi.org/10.1016/S0896-6273\(02\)00569-X](https://doi.org/10.1016/S0896-6273(02)00569-X).
46. Westman E, Aguilar C, Muehlboeck J-S, Simmons A. Regional magnetic resonance imaging measures for multivariate analysis in Alzheimer's disease and Mild Cognitive Impairment. *Brain Topogr*. 2012;26(1):9–23. <https://doi.org/10.1007/s10548-012-0246-x>.
47. Kalavathi P, Senthamilselvi M, Prasath VBS. Review of computational methods on brain symmetric and asymmetric analysis from neuroimaging techniques. *Technologies*. 2017. <https://doi.org/10.3390/technologies5020016>.
48. Roe JM, Vidal-Piñeiro D, Sørensen Ø, Brandmaier AM, Düzel S, Gonzalez HA, Kievit RA, Knights E, Kühn S, Lindenberger U, Mowinckel AM, Nyberg L, Park DC, Pudas S, Rundle MM, Walhovd KB, Fjell AM, Westerhausen R. Asymmetric thinning of the cerebral cortex across the adult lifespan is accelerated in Alzheimer's disease. *Nat Comm*. 2021;12(1). <https://doi.org/10.1038/s41467-021-21057-y>.
49. Wu X, Wu Y, Geng Z, Zhou S, Wei L, Ji G-J, Tian Y, Wang K. Asymmetric differences in the gray matter volume and functional connections of the amygdala are associated with clinical manifestations of Alzheimer's disease. *Front Neurosci*. 2020;14:602. <https://doi.org/10.3389/fnins.2020.00602>.
50. Low A, Mak E, Malpetti M, Chouliaras L, Nicastro N, Su L, Holland N, Rittman T, Rodríguez PV, Passamonti L, Bevan-Jones WR, Jones PS, Rowe JB, O'Brien JT. Asymmetrical atrophy of thalamic subnuclei in Alzheimer's disease and amyloid-positive Mild Cognitive Impairment is associated with key clinical features. *Alzheimers Dement (Amst)*. 2019;11(1):690–9. <https://doi.org/10.1016/j.dadm.2019.08.001>.
51. Strobl C, Boulesteix A-L, Zeileis A, Hothorn T. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*. 2007;8(1):1–21. <https://doi.org/10.1186/1471-2105-8-25>.
52. Moćkus J. On bayesian methods for seeking the extremum. In: *Proceedings of the Optimization Techniques IFIP Technical Conference, 1975;400–4*. Springer, Berlin. https://doi.org/10.1007/3-540-07165-2_55
53. Head T, Kumar M, Nahrstaedt H, Louppe G, Shcherbatyi I. scikit-optimize/scikit-optimize. Zenodo. 2020. <https://doi.org/10.5281/zenodo.4014775>.
54. McKay MD, Beckman RJ, Conover WJ. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*. 1979;21(2):239–45. <https://doi.org/10.2307/1268522>.
55. Bloch L, Friedrich CM. Using bayesian optimization to effectively tune random forest and XGBoost hyperparameters for early Alzheimer's disease diagnosis. In: Ye J, O'Grady MJ, Civitarese G, Yordanova K (eds) *Wireless Mobile Communication and Healthcare, 2021;285–99*. Springer, Basel, Switzerland. https://doi.org/10.1007/978-3-030-70569-5_18.
56. Refaeilzadeh P, Tang L, Liu H. Cross-validation. In: Liu L, Özsu MT (eds) *Encyclopedia of Database Systems*, pp. 532–8. Springer, Boston, Massachusetts, United States 2009. https://doi.org/10.1007/978-0-387-39940-9_565.
57. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*. 2011;12:2825–30.
58. Santos MS, Soares JP, Abreu PH, Araujo H, Santos J. Cross-validation for imbalanced datasets: Avoiding overoptimistic and overfitting approaches. *IEEE Comput Intell Mag*. 2018;13(4):59–76. <https://doi.org/10.1109/MCI.2018.2866730>.
59. Chen T, He T, Benesty M, Khotilovich V, Tang Y, Cho H, Chen K, Mitchell R, Cano I, Zhou T, Li M, Xie J, Lin M, Geng Y, Li Y. xgboost: Extreme Gradient Boosting. (2019). Manual of R package v0.82.1 <https://CRAN.R-project.org/package=xgboost>, Accessed: 2021-10-10.
60. Efron B, Tibshirani R. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Stat Sci*. 1986;1(1):54–75. <https://doi.org/10.1214/ss/1177013815>.
61. Rosasco L, Vito ED, Caponnetto A, Piana M, Verri A. Are loss functions all the same? *Neural Comput*. 2004;16(5):1063–76. <https://doi.org/10.1162/089976604773135104>.
62. Breiman L, Friedman J, Stone CJ, Olshen RA. *Classification and Regression Trees*, 1st edn. CRC press, New York, New York, US 1984. <https://doi.org/10.1201/9781315139470>.
63. Cox DR. The regression analysis of binary sequences. *J R Stat Soc Series B Stat Methodol*. 1958;20(2):215–32. <https://doi.org/10.1111/j.2517-6161.1958.tb00292.x>.
64. Shapley LS. A value for n-person games. In: Kuhn HW, Tucker AW (eds) *Contributions to the Theory of Games (AM-28)* 1953;2:307–18. Princeton University Press, Princeton, New Jersey, US. <https://doi.org/10.1515/9781400881970-018>.
65. Štrumbelj E, Kononenko I. Explaining prediction models and individual predictions with feature contributions. *Knowl Inf Syst*. 2013;41(3):647–65. <https://doi.org/10.1007/s10115-013-0679-x>.
66. Ribeiro M, Singh S, Guestrin C. Why should I trust you?: Explaining the predictions of any classifier. In: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations (NAACL 2016): 12-17 June 2016; San Diego, California, US, 2016;97–101*. Association for Computational Linguistics, San Diego, US. <https://doi.org/10.18653/v1/n16-3020>.
67. Lundberg SM, Erion GG, Lee S. Consistent individualized feature attribution for tree ensembles. *Computing Research Repository* 2018. [arxiv:1802.03888](https://arxiv.org/abs/1802.03888), Accessed: 2021-10-10.
68. Kumar IE, Venkatasubramanian S, Scheidegger C, Friedler S. Problems with Shapley-value-based explanations as feature importance measures. In: *Proceedings of the International Conference on Machine Learning (ICML), 2020;5491–500*. PMLR. <http://proceedings.mlr.press/v119/kumar20e/kumar20e.pdf>.
69. Pekala K, Woznica K, Biecek P. Triplot: Model agnostic measures and visualisations for variable importance in predictive

- models that take into account the hierarchical correlation structure. CoRR abs/2104.03403. 2021.
70. Frisoni GB, Fox NC, Jack CR, Scheltens P, Thompson PM. The clinical use of structural MRI in Alzheimer disease. *Nat Rev Neurol*. 2010;6(2):67–77. <https://doi.org/10.1038/nrneuro.2009.215>.
 71. Mueller SG, Schuff N, Yaffe K, Madison C, Miller B, Weiner MW. Hippocampal atrophy patterns in Mild Cognitive Impairment and Alzheimer's disease. *Hum Brain Mapp*. 2010;31(9):1339–47. <https://doi.org/10.1002/hbm.20934>.
 72. deToledo-Morrell L, Stoub TR, Bulgakova M, Wilson RS, Bennett DA, Leurgans S, Wu J, Turner DA. MRI-derived entorhinal volume is a good predictor of conversion from MCI to AD. *Neurobiol Aging*. 2004;25(9):1197–203. <https://doi.org/10.1016/j.neurobiolaging.2003.12.007>.
 73. Poulin SP, Dautoff R, Morris JC, Barrett LF, Dickerson BC. Amygdala atrophy is prominent in early Alzheimer's disease and relates to symptom severity. *Psychiatry Res Neuroimaging*. 2011;194(1):7–13. <https://doi.org/10.1016/j.psychresns.2011.06.014>.
 74. Yang H, Xu H, Li Q, Jin Y, Jiang W, Wang J, Wu Y, Li W, Yang C, Li X, Xiao S, Shi F, Wang T. Study of brain morphology change in Alzheimer's disease and amnesic Mild cognitive impairment compared with normal controls. *General Psychiatry*. 2019. <https://doi.org/10.1136/gpsych-2018-100005>.
 75. Herzog NJ, Magoulas GD. Brain asymmetry detection and machine learning classification for diagnosis of early dementia. *Sensors*. 2021. <https://doi.org/10.3390/s21030778>.
 76. Foundas AL, Leonard CM, Mahoney SM, Agee OF, Heilman KM. Atrophy of the hippocampus, parietal cortex, and insula in alzheimer's disease: a volumetric magnetic resonance imaging study. *Neuropsychiatry Neuropsychol Behav Neurol*. 1997;10(2):81–9.
 77. Greene SJ, Killiany RJ. Subregions of the inferior parietal lobule are affected in the progression to Alzheimer's disease. *Neurobiol Aging*. 2010;31(8):1304–11. <https://doi.org/10.1016/j.neurobiolaging.2010.04.026>.
 78. Yao Z, Zhang Y, Lin L, Zhou Y, Xu C, Jiang T. the Alzheimer's Disease Neuroimaging Initiative: Abnormal cortical networks in Mild Cognitive Impairment and Alzheimer's disease. *PLoS Comput Biol*. 2010;6(11):1–11. <https://doi.org/10.1371/journal.pcbi.1001006>.
 79. Scheff SW, Price DA, Schmitt FA, Scheff MA, Mufson EJ. Synaptic loss in the inferior temporal gyrus in Mild Cognitive Impairment and Alzheimer's disease. *J Alzheimers Dis*. 2011;24(3):547–57. <https://doi.org/10.3233/JAD-2011-101782>.
 80. Visser PJ, Verhey FRJ, Hofman PAM, Scheltens P, Jolles J. Medial temporal lobe atrophy predicts Alzheimer's disease in patients with minor cognitive impairment. *J Neurol Neurosurg Psychiatry*. 2002;72(4):491–7. <https://doi.org/10.1136/jnnp.72.4.491>.
 81. Yokoi T, Watanabe H, Yamaguchi H, Bagarinao E, Masuda M, Imai K, Ogura A, Ohdake R, Kawabata K, Hara K, Riku Y, Ishigaki S, Katsuno M, Miyao S, Kato K, Naganawa S, Harada R, Okamura N, Yanai K, Yoshida M, Sobue G. Involvement of the precuneus / posterior cingulate cortex is significant for the development of Alzheimer's disease: A PET (THK5351, PiB) and resting fMRI study. *Front Aging Neurosci*. 2018;10. <https://doi.org/10.3389/fnagi.2018.00304>.
 82. Tabatabaei-Jafari H, Shaw ME, Cherbuin N. Cerebral atrophy in Mild Cognitive Impairment: A systematic review with meta-analysis. *Alzheimers Dement (Amst)*. 2015;1(4):487–504. <https://doi.org/10.1016/j.dadm.2015.11.002>.
 83. Zhang Y, Schuff N, Camacho M, Chao LL, Fletcher TP, Yaffe K, Woolley SC, Madison C, Rosen HJ, Miller BL, Weiner MW. MRI markers for Mild Cognitive Impairment: Comparisons between white matter integrity and gray matter volume measurements. *PLoS ONE*. 2013;8(6):1–10. <https://doi.org/10.1371/journal.pone.0066367>.
 84. Ledig C, Schuh A, Guerrero R, Heckemann RA, Rueckert D. Structural brain imaging in Alzheimer's disease and mild cognitive impairment: biomarker analysis and shared morphometry database. *Sci Rep*. 2018. <https://doi.org/10.1038/s41598-018-29295-9>.
 85. Thompson PM, Hayashi KM, de Zubicaray GI, Janke AL, Rose SE, Semple J, Hong MS, Herman DH, Gravano D, Dreddell DM, Toga AW. Mapping hippocampal and ventricular change in Alzheimer disease. *Neuroimage*. 2004;22(4):1754–66. <https://doi.org/10.1016/j.neuroimage.2004.03.040>.
 86. Jack CR, Shiung MM, Gunter JL, O'Brien PC, Weigand SD, Knopman DS, Boeve BF, Ivnik RJ, Smith GE, Cha RH, Tangalos EG, Petersen RC. Comparison of different MRI brain atrophy rate measures with clinical disease progression in AD. *Neurology*. 2004;62(4):591–600. <https://doi.org/10.1212/01.WNL.0000110315.26026.EF>.
 87. Corder E, Saunders A, Strittmatter W, Schmechel D, Gaskell P, Small G, Roses A, Haines J, Pericak-Vance M. Gene dose of Apolipoprotein E type 4 allele and the risk of Alzheimer's disease in late onset families. *Science*. 1993;261(5123):921–3. <https://doi.org/10.1126/science.8346443>.
 88. Strittmatter WJ, Saunders AM, Schmechel D, Pericak-Vance M, Enghild J, Salvesen GS, Roses AD. Apolipoprotein E: High-avidity binding to beta-amyloid and increased frequency of type 4 allele in late-onset familial Alzheimer disease. *Proc Natl Acad Sci USA*. 1993;90(5):1977–81. <https://doi.org/10.1073/pnas.90.5.1977>.
 89. Elias-Sonnenschein LS, Viechtbauer W, Ramakers IHGB, Verhey FRJ, Visser PJ. Predictive value of APOE-ε4 allele for progression from MCI to AD-type dementia: A meta-analysis. *J Neurol Neurosurg Psychiatry*. 2011;82(10):1149–56. <https://doi.org/10.1136/jnnp.2010.231555>.
 90. Minkova L, Habich A, Peter J, Kaller CP, Eickhoff SB, Klöppel S. Gray matter asymmetries in aging and neurodegeneration: A review and meta-analysis. *Hum Brain Mapp*. 2017;38(12):5890–904. <https://doi.org/10.1002/hbm.23772>.
 91. Wachinger C, Salat DH, Weiner M, Reuter M. for the Alzheimer's Disease Neuroimaging Initiative: Whole-brain analysis reveals increased neuroanatomical asymmetries in dementia for hippocampus and amygdala. *Brain*. 2016;139(12):3253–66. <https://doi.org/10.1093/brain/aww243>.
 92. Wachinger C, Golland P, Kremen W, Fischl B, Reuter M. BrainPrint: A discriminative characterization of brain morphology. *Neuroimage*. 2015;109:232–48. <https://doi.org/10.1016/j.neuroimage.2015.01.032>.
 93. Friedman M. The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *J Am Stat Ass*. 1937;32(200):675–701. <https://doi.org/10.1080/01621459.1937.10503522>.
 94. Kendall MG. A new measure of rank correlation. *Biometrika*. 1938;30(1–2):81–93. <https://doi.org/10.1093/biomet/30.1-2.81>.
 95. Merkel D. Docker: Lightweight Linux containers for consistent development and deployment. *Linux J*. 2014;2014(239):2.
 96. Molnar C, Casalicchio G, Bischl B. Quantifying model complexity via functional decomposition for better post-hoc interpretability. In: Cellier P, Driessens K, editors. *Machine Learning and Knowledge Discovery in Databases*. Berlin: Springer; 2020. p. 93–204. https://doi.org/10.1007/978-3-030-43823-4_17.
 97. Lukasová A. Hierarchical agglomerative clustering procedure. *Pattern Recognit*. 1979;11(5):365–81. [https://doi.org/10.1016/0031-3203\(79\)90049-9](https://doi.org/10.1016/0031-3203(79)90049-9).
 98. Baniecki H, Kretowicz W, Piatyszek P, Wisniewski J, Biecek P. dalex: Responsible machine learning with interactive explainability and fairness in Python. *J Mach Learn Res*. 2021;22(214):1–7.

99. Yao Z, Hu B, Liang C, Zhao L, Jackson M. the Alzheimer's Disease Neuroimaging Initiative: A longitudinal study of atrophy in amnesic Mild Cognitive Impairment and normal aging revealed by cortical thickness. PLoS ONE. 2012;7(11):1–11. <https://doi.org/10.1371/journal.pone.0048973>.
100. Beekley DL, Ramos EM, van Belle G, Deitrich W, Clark AD, Jacka ME, Kukull WA. The National Alzheimer's Coordinating Center (NACC) database: An Alzheimer Disease database. Alzheimer Dis. Assoc. Disord. 2004;18(4), 270–7. https://journals.lww.com/alzheimerjournal/Abstract/2004/10000/The_National_Alzheimer_s_Coordinating_Center.21.aspx, Accessed: 2021-10-10.
101. Snoek J, Larochelle H, Adams RP. Practical bayesian optimization of machine learning algorithms. In: Pereira F, Burges CJC, Bottou L, Weinberger KQ (eds) Advances in Neural Information Processing Systems (NIPS), 2012;25, 2951–9. Curran Associates, Inc., New York, New York, US. <https://proceedings.neurips.cc/paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf> Accessed 2021-10-10.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.