



Alzheimer's Disease Classification using ML Pipeline on Fast Fourier Transformed EEG Data



Tyler Yoshihara

Pomona College Department of Neuroscience - NEUR182 - Machine Learning with Neural Signals

Abstract

Alzheimer's Disease (AD) is the most common neurodegenerative disease. It is typically late onset and can develop substantially before diagnosable symptoms appear. Electroencephalogram (EEG) could potentially serve as a noninvasive diagnostic tool for AD. Machine learning can be helpful in making inferences about changes in frequency bands in EEG data and how these changes relate to neural function. The EEG data was sourced from 2014 paper titled *Alzheimer's disease patients classification through EEG signals processing* by Fison et al. There were patients with AD, mild cognitive impairment (MCI), and healthy controls. The data was already preprocessed using a fast fourier transform (FFT) to take the data from the time domain to the frequency domain. There were differing levels of effectiveness in terms of classification but generally, Fisher's discriminant analysis (FDA), relevance vector machine, and random forest approaches were most successful. Due to inconsistent feature importances in different models, conclusions about important frequency bands for classification were not able to be made at this time. Similarly, different frequencies were not able to be localized to different regions of the brain. Further research is necessary to develop more interpretable models for classification.

Introduction

Alzheimer's disease is an irreversible, progressive neurodegenerative disorder marked by memory issues associated with dementia, language problems, and erratic behavior. While diagnostic techniques such as identifying the biomarkers in the brain have been proposed as an indicator of a patient's proclivity to developing Alzheimer's, inaccuracy is a problem that plagues a lot of these approaches (Tarnanas, Tsolaki, Wiederhold, Wiederhold, & Tsolaki, 2015). Currently, invasive, posthumous brain examination is the only guaranteed diagnostic tool. As a result, the application of machine learning to EEG data could be an extremely useful and cost-effective method to non-invasively screen for Alzheimer's disease and could potentially serve as an incredibly beneficial and life-saving medical protocol (Ding et al., 2018). Often Alzheimer's symptoms are associated with later stages of the disease, which is why early screening could allow for intervention prior to neurodegeneration running its course. The classification technique for Alzheimer's disease necessitates cross-subject analysis, as intra-subject classification would not provide any useable results. The EEG data used for the body of this study comes from a 2018 study by Fison et al., utilizing classification techniques derived in an earlier 2014 study by the same group (Fison et al., 2014). Our hope in expanding upon Fison et al. 2018's work is to improve upon their classification techniques and provide interpretable results that can be localized to specific frequency bands and channels for diagnostic purposes (Fison et al., 2018). Fison et al. employed decision trees (DT), support vector machines (SVC) and rule-based classifiers, ultimately citing DTs as their most effective algorithm in their 2018 paper. In selecting our models, the primary goal was to maximize sensitivity and specificity, while retaining parsimony and being highly cognizant of overfitting and decreases in accuracy due to features derived from noise that is not consistent across cases (Hebart & Baker, 2018). Our classification algorithms include Relevance Vector Classifier (RVC), Ridge Regularized Linear Regression (RLR - L2), Fisher's Discriminant Analysis (FDA), and Random Forest (RF) (Liaw et al., 2002). Four primary EEG frequency bands were examined for salience in our interpretation of feature importance: delta (0.5-4 Hz), theta (4-7 Hz), alpha (8-13 Hz), beta (13-30 Hz). It was hypothesized that alpha bands would be downregulated across channels, while beta and delta frequency would increase for AD patients relative to healthy control. MCI patients were expected to have less degradation of alpha frequency compared to AD subjects.

Existing Data Set

The existing data, taken from Fison et al. 2018's study, is pre-processed EEG recording from 86 participants with AD or Mild Cognitive Impairment (MCI) and 23 Healthy Controls (HC). The data was derived from a resting-eyes closed (EC) state for 300 seconds, with the central three-minutes of data (60 sec to 240 sec) being the focus of analysis. An 19 electrode array with a sampling rate of 256 Hz was used. Fison et al. preprocessed the data utilizing two techniques, one of which will be focused on in this follow-up study. Fast Fourier Transform (FFT) was conducted on the central three-minutes of data, which was divided into six, 30 second epochs. 16 Fourier coefficients were extracted from each epoch, which was presented to us in the form of a csv file with cases represented and labeled on the y-axis and features in the form of fourier coefficients on the x-axis (304 columns).

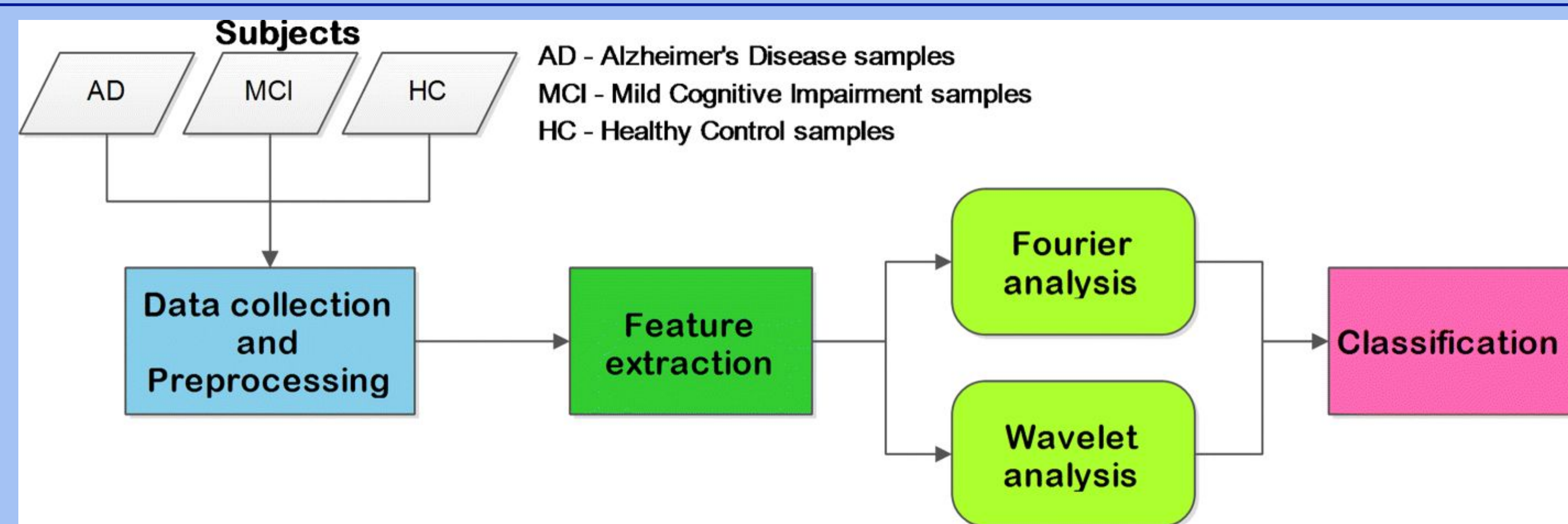


Figure 1: Methodology taken from Fison et al. 2018 - (Data was presented in the a feature extracted, pre-processed state, ready for classification.)

Methods

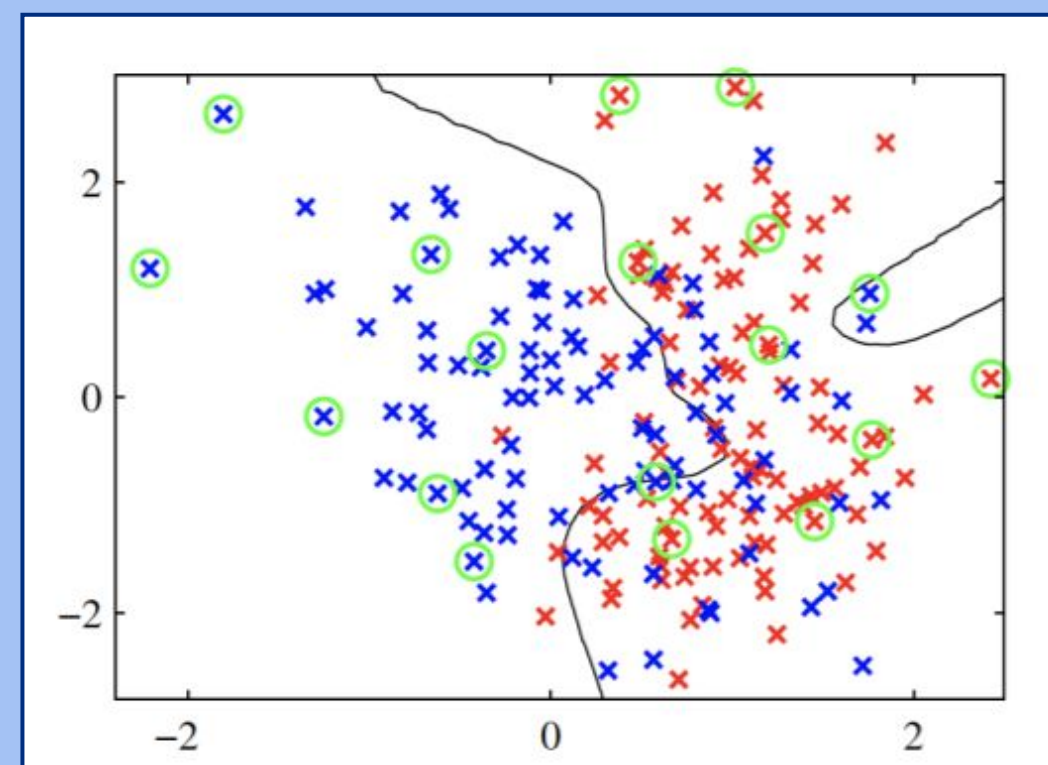


Figure 2 : Relevance Vector Classifier - drawn decision boundary and circled relevance vectors

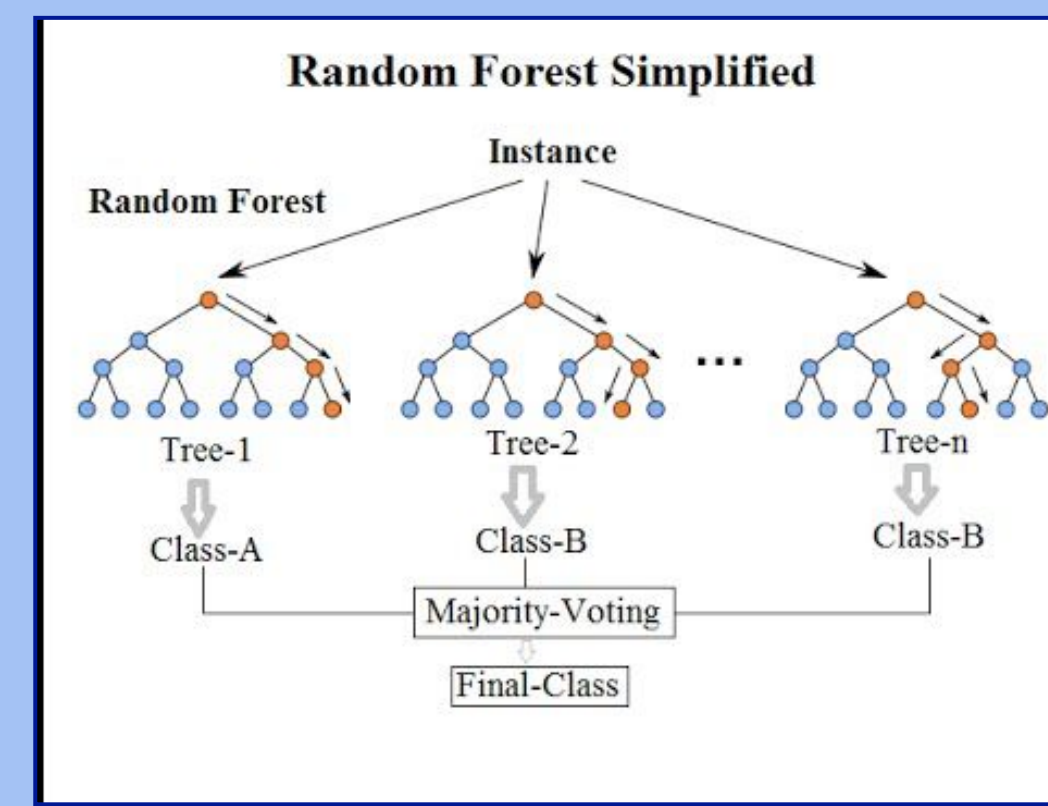


Figure 3 : Random Forest Classifier - Disjoint set of decision trees that are combined to form an ensemble method for classification

Machine Learning Approaches:

- Relevance Vector Classifier (RVC)** - Bayesian sparse kernel method that builds a separating hyperplane using probabilistic measures to maximizes the minimum distance between the data of different classes (Bishop, 2006).
 - Relies on significantly less basis functions than an SVM and, therefore, provides a much more parsimonious solution
- Random Forest (RF)** - Ensemble method that builds and combines decision trees while searching for the best features among random subsets of features (Liaw & Wiener, 2001).
 - Fison et al. had their best results with a simple decision tree. Because decision trees form the framework for random forests, it is expected that this will also be an effective method.
- Fisher's Discriminant Analysis (FDA)** - Using single value decomposition to transform feature vectors into a space that maximizes separability.

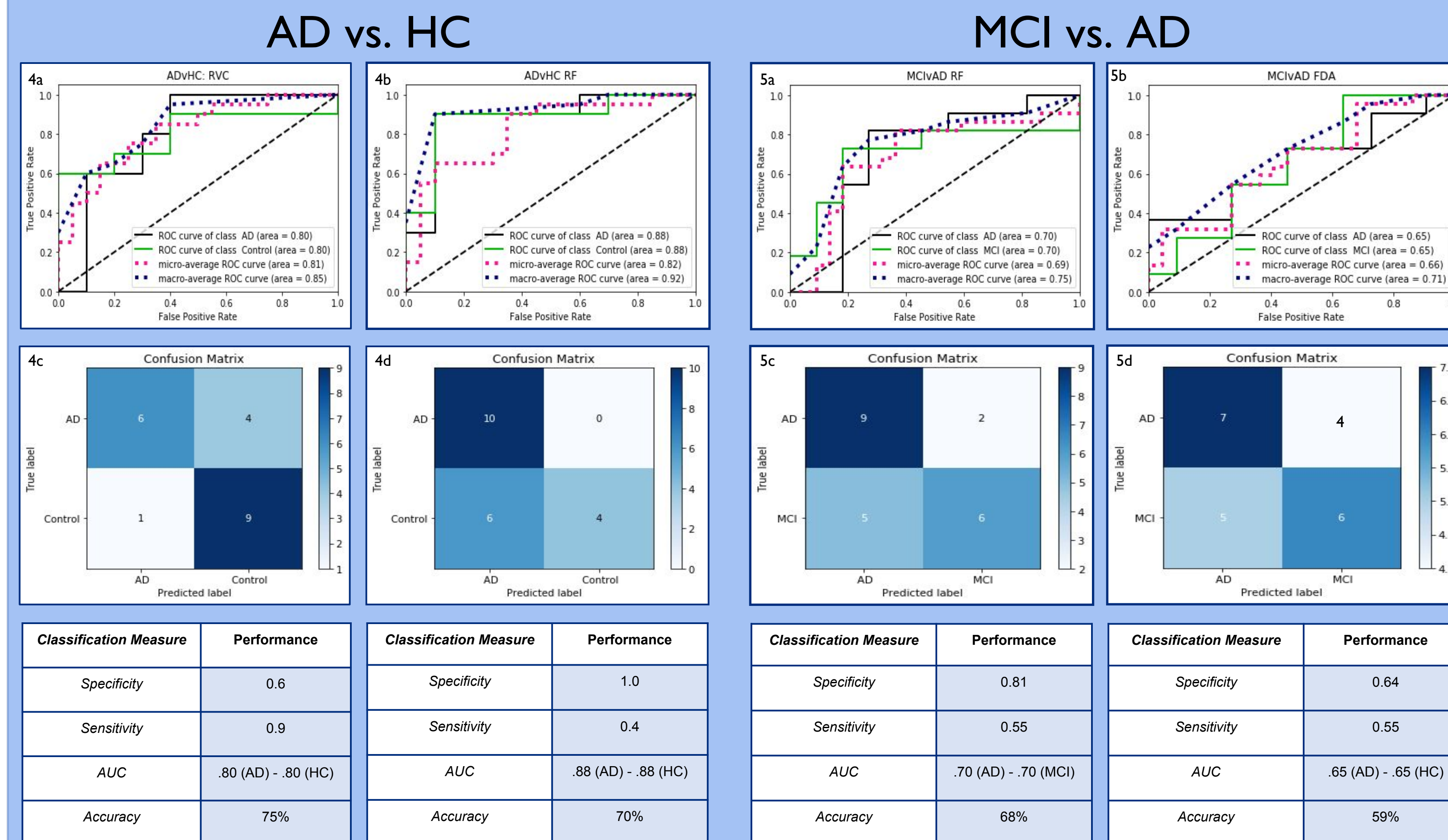
Model Accuracy:

- Accuracy, ROC, Sensitivity, Specificity

Interpretation:

- Cross-model feature importance correlation
- Frequency band analysis
- Spatial analysis by frequency and channel

Results



Figures 4a,b,c,d: AD vs. HC Classification and Accuracy Measures - RVC, RLR: L2, RF & FDA were utilized and the results of RVC (a, c) & RF (c, d) are presented above.

HC vs. MCI

Classification Measure	Performance
Specificity	0.8
Sensitivity	0.71
AUC	.60 (HC) - .60 (MCI)
Accuracy	75%

RF

RVC

Figures 6a,b: HC vs. MCI Accuracy Measures - RVC, RLR, L2 & FDA were utilized and the results of RF (a) and RVC (b) are presented to the left.

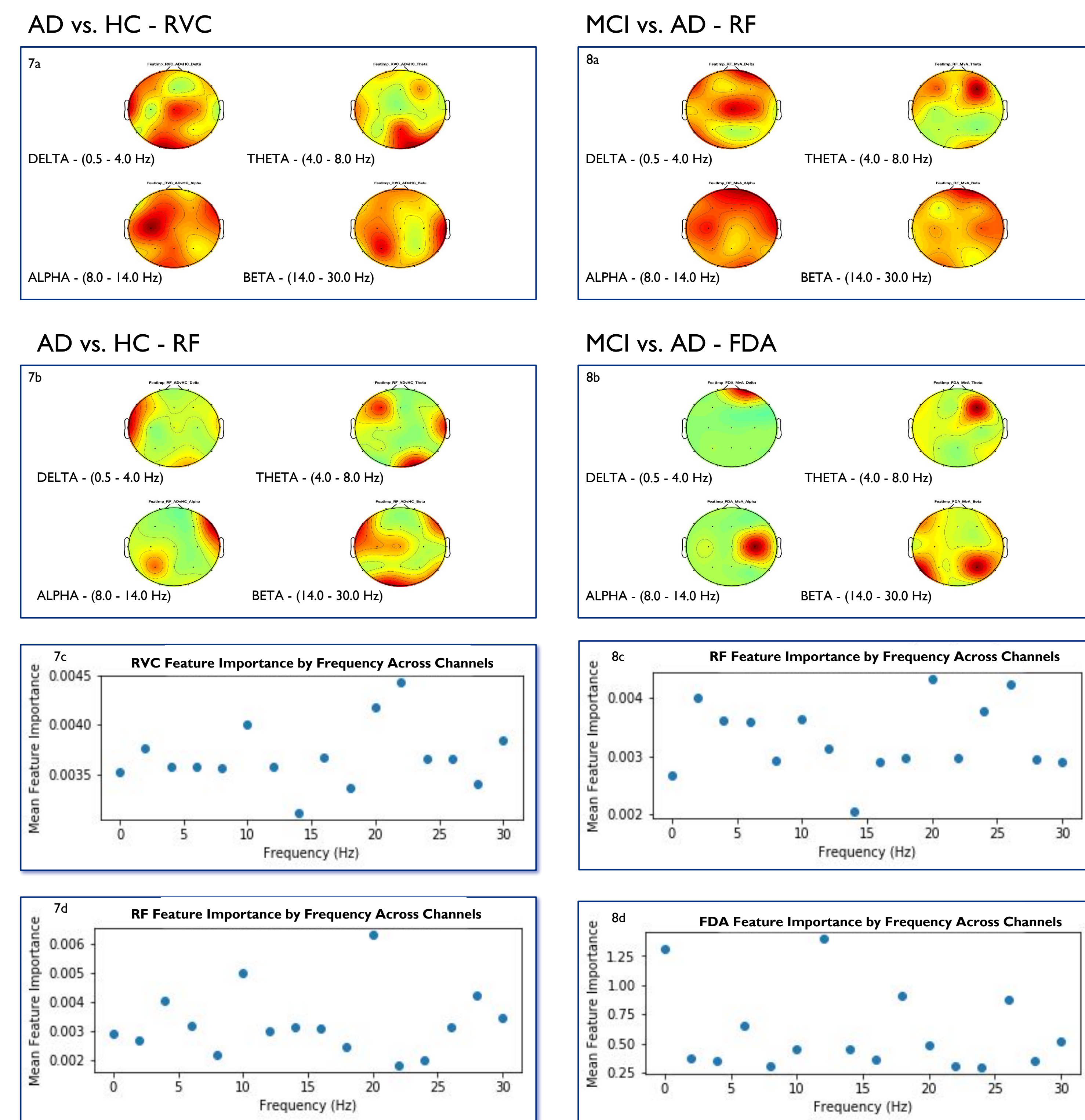


Fig 7a,b,c,d: Heatmaps and Feature Importance for AD vs. HC - Relative feature importances did not demonstrate any consistency across channels and frequencies

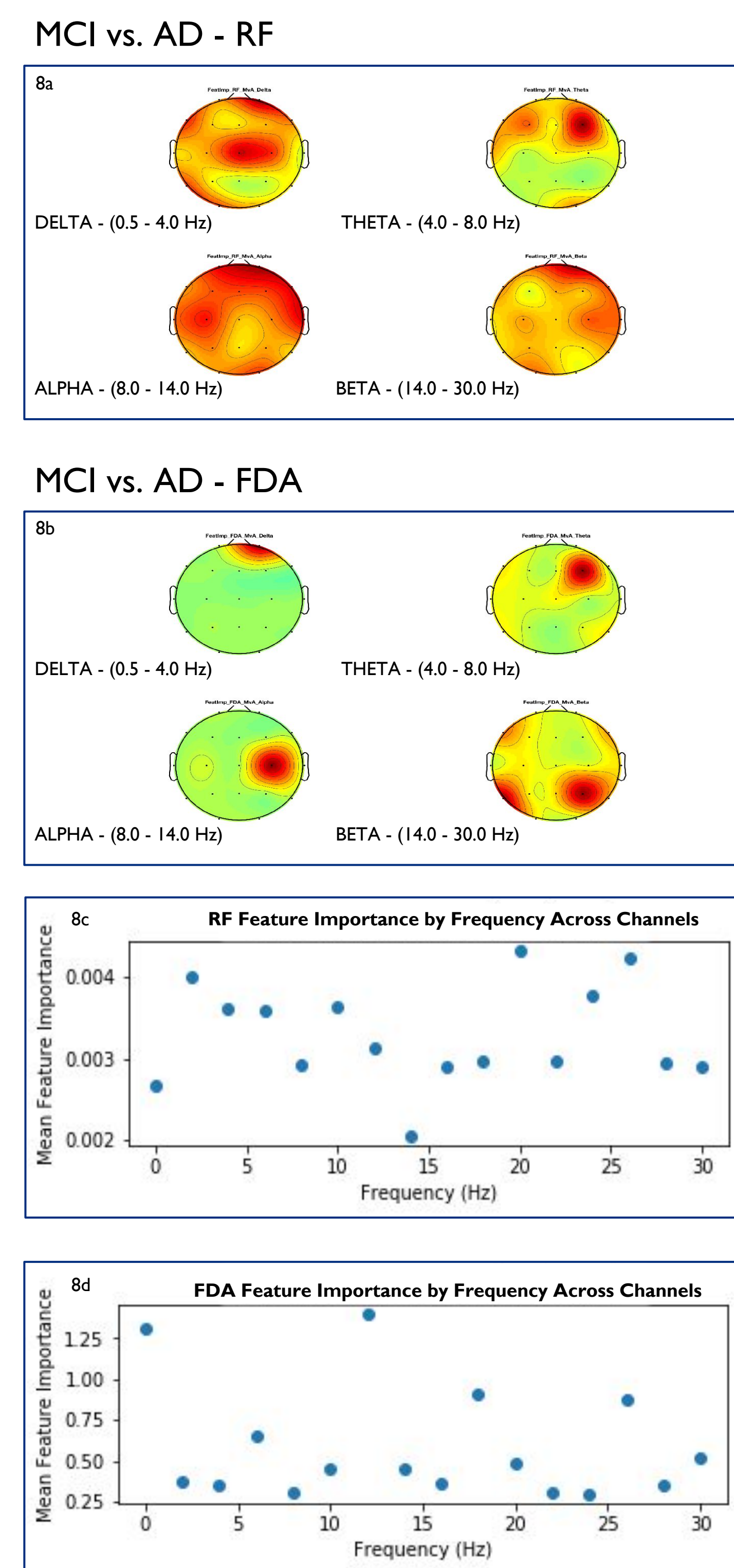


Fig 8a,b,c,d: Heatmaps and Feature Importance for MCI vs. AD - Relative feature importances did not demonstrate any consistency across channels and frequencies

Conclusion

- RVC was most effective for AD vs. HC
- RF was most effective for MCI vs AD and HC vs MCI
- PCA was not utilized for RLR or FDA as there was minimal multicollinearity between features
- Correlations between feature importance for different algorithms within dataset was low
 - Conclusions not able to be drawn about consistency of feature importance between models
- Correlations between feature importance for frequency bands across channels was inconsistent
 - Heat maps for different models within data was not consistent in terms of spatial orientation
 - Different frequencies were not localized to specific regions and were not overlapping between algorithms
- Future research is necessary in order to draw more interpretable conclusions pertaining to locality and frequency band specificity when it comes to Alzheimer's classification.

Acknowledgments

Thank you to Professor Spezio for guidance in this project and thank you to Fison et al. for providing the preprocessed data and original study that inspired this work. Lastly, thank you to Isabelle and Fernanda for collaborating with us in the early stages of this project.

References

- Bishop, C. (2006). Pattern recognition and machine learning. Springer.
- Ding, Y., Saha, J. B., Kowalewski, M. G., Trivedi, R., Harnish, R., Jenkins, N. W., ... Fison, G. (2019). A Deep Learning Model to Predict a Diagnosis of Alzheimer's Disease by Using 18F-FDG PET of the Brain. *Neurology*, 290(2), 456-464. <https://doi.org/10.1213/01.neuro.0000595955>
- Fison, G., Weischek, E., Calhoun, A., Felis, G., Benitez, P., De Salvo, S., ... De Cola, M. C. (2018). A Combining EEG signal processing with supervised methods for Alzheimer's patients classification. *BMC Medical Informatics and Decision Making*, 18(1), 35. <https://doi.org/10.1186/s12911-018-0455-5>
- Fison, G., Weischek, E., Felis, G., Benitez, P., De Salvo, S., Benatti, P., & De Cola, M. C. (2014, December 10). *Alzheimer's disease patients classification through EEG signals processing*. <https://doi.org/10.1109/CICWD412014.2014.7006655>
- Fisher, M. A., & Malar, C. I. (2010). Determining multivariate decoding for the study of brain function. *NeuroImage*, 180, 4-18. <https://doi.org/10.1016/j.neuroimage.2017.08.040>
- Liaw, A., & Wiener, M. (2001). Classification and Regression by Random Forests. *Forest*, 23.
- Tarnanas, I., Tsolaki, A., Wiederhold, M., Wiederhold, B., & Tsolaki, M. (2015). Five-year biomarker progression variability for Alzheimer's disease dementia prediction: Can a complex instrumental activities of daily living marker fill in the gaps? *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring*, 1(4), 521-532. <https://doi.org/10.1016/j.dadm.2015.10.005>