













Alzheimer Disease Detection from Raman Spectroscopy of the Cerebrospinal Fluid via Topological Machine Learning

Francesco Conti^{1,2}, Martina Banchelli³, Valentina Bessi⁴, Cristina Cecchi⁵, Fabrizio Chiti⁵, Sara Colantonio¹, Cristiano D'Andrea³, Marella de Angelis³, Davide Moroni¹, Benedetta Nacmias^{4,6}, Maria Antonietta Pascali¹, Sandro Sorbi^{4,6} and Paolo Matteini³

¹ Institute of Information Science and Technologies “A. Faedo”,
National Research Council of Italy, Pisa (IT)
`{Name.Surname}@isti.cnr.it`

² Department of Mathematics, University of Pisa, Pisa (IT)

³ Institute of Applied Physics “N. Carrara”, National Research Council of Italy,
Sesto Fiorentino (IT)

⁴ Department of Neuroscience, Psychology, Drug Research and Child Health,
University of Florence, Florence (IT)

⁵ Department of Clinical and Experimental Biomedical Sciences “Mario Serio”,
University of Florence, Florence (IT)

⁶ IRCCS Fondazione Don Carlo Gnocchi, Florence (IT)

Abstract. The cerebrospinal fluid (CSF) of 19 subjects who received a clinical diagnosis of Alzheimer’s disease (AD) as well as of 5 pathological controls have been collected and analysed by Raman spectroscopy (RS). We investigated whether the raw and preprocessed Raman spectra could be used to distinguish AD from controls. First, we applied standard Machine Learning (ML) methods obtaining unsatisfactory results. Then, we applied ML to a set of topological descriptors extracted from raw spectra, achieving a very good classification accuracy ($> 87\%$). Although our results are preliminary, they indicate that RS and topological analysis together may provide an effective combination to confirm or disprove a clinical diagnosis of AD. The next steps will include enlarging the dataset of CSF samples to validate the proposed method better and, possibly, to understand if topological data analysis could support the characterization of AD subtypes.

Keywords: Ensembling · bagging · machine learning · deep learning · image classification · convolutional neural networks.

1 Introduction

Alzheimer’s disease (AD) affects tens of millions of people worldwide, being the most common neurodegenerative disease. Due to the population aging, the number of people affected by AD and other forms of dementia is expected

to reach about 152 million by 2050 (World Alzheimer Report 2021 provided by Alzheimer’s Disease International, McGill University <https://www.alzint.org/resource/world-alzheimer-report-2021/>). At present, the clinical diagnosis of AD requires a series of neurological examinations (National Institute of Aging – Alzheimer’s Association criteria), while the definitive diagnosis is possible only after the patient’s death and brain tissue analysis. Therefore, there is a need to improve the accuracy of clinical diagnosis with innovative, cost-effective and specific approaches. Raman spectroscopy (RS) represents a fast, efficient, non-invasive diagnostic tool [9], and the high-precision detection of RS is expected to reduce or replace other AD diagnostic tests. Recently, Raman-based techniques demonstrated significant potential in identifying AD by detecting specific biomarkers in body fluids [13]. Given the increasing number of RS studies, a systematic evaluation of the accuracy of RS in the diagnosis of AD was already performed, showing that RS is an effective and accurate tool for diagnosing AD, though it still cannot rule out the possibility of misdiagnosis [19]. Recently, Raman spectroscopy of tissue samples has been coupled with topological machine learning to support the grading of bone cancer [6], showing the feasibility of a topological approach for multi-label classification. The detection of CSF biomarkers is one of the diagnostic criteria for AD [2] because CSF is more sensitive than blood or other biofluids in the diagnosis of AD. Therefore, RS can be used as an effective tool to analyze CSF samples, as shown previously [14,12].

Here we propose a novel method based on the collection of the vibrational Raman fingerprint of the proteomic content of cerebrospinal fluid (CSF) and on the topological machine learning analysis of the Raman spectra in order to support the AD diagnosis. The achieved results encourage to keep on investigating topological machine learning tools, not only to establish more safely the proposed methodology by enlarging the experimentation but also to understand if looking at Raman spectra of CSF with the topological lens could also help to characterize AD subtypes.

2 Population study and Data Acquisition

The study population is made of 24 patients, enrolled in the framework of the Bando Salute 2018 PRAMA project (“Proteomics, RAdiomics & Machine learning-integrated strategy for precision medicine for Alzheimer’s”), co-funded by the Tuscany Region, with the approval of the Institutional Ethics Committee of the Careggi University Hospital Area Vasta Centro (ref. number 17918_bio). All of them showed pathological symptoms: the majority of them, 19 subjects, have been diagnosed with AD, while the others have been considered as controls (noAD), even if diagnosed with other neurological conditions: one with vascular dementia, three with hydrocephalus and one with multiple sclerosis.

The CSF samples were collected by lumbar puncture, then immediately centrifuged at 200*g* for 1 min, 20 °C and stored at –80 °C until analysis [15,17]. On the day of analysis, CSF samples were thawed and centrifuged again at 4000*g*

for 10 min at 4 °C. The pellet was separated from the supernatant and further used for the analyses. A 2 ul drop of the pellet was deposited onto a gold mirror support (ME1S-M01; Thorlabs, Inc., Newton, NJ), followed by air drying for 30 minutes and acquisition of Raman spectra from the outer ring of the dried drop. A set of five Raman spectra have been collected for each drop-casted sample by using a micro-Raman spectrometer (Horiba, France) in back-scattering configuration, equipped with a laser excitation source tuned at 785 nm (40 mW power, 20 second integration time, 10 accumulations) and a Peltier cooled CCD detector.

Finally, a set of five Raman spectra have been collected for each biological sample. In some cases, the same procedure has been replicated two or three times; it resulted in a dataset of 30 acquisitions of RS: 22 belonging to the AD class and 8 to the noAD class.

3 Methods

After the Raman spectra are acquired, the data enter the following pipeline to return the final predictive model with classification accuracy. For each patient, the average of the five acquisitions of the raw Raman spectrum is computed. Next, the following transformations are applied to the RS: Fourier transform, Welch transform and autocorrelation. We applied the pipeline individually on the original spectra and on each of the transformations listed above. These computations were performed using the Python package SciPy [18].

The spectra enter the pipeline of Topological Machine Learning (TML). For more detailed information on the pipeline, refer to [7]. The pipeline performs a lower star filtration to extract the Persistence Diagrams (PDs). Since the data is a 1D spectrum, the only non-trivial homology group is H_0 . The PD is vectorized using the following vectorization methods: Persistence Image [1] with parameters $\sigma \in \{0.1, 1, 10\}$, $n \in \{5, 10, 25\}$, Persistence Landscape [4], Persistence Silhouette [5] and Betti curve [16] all with parameters $n \in \{25, 50, 75, 100\}$. Finally, these vectors enter one of the following Machine Learning (ML) classifiers: Support Vector Classifier [8], Random Forest Classifier [3] and Ridge Classifier [11]. The validation scheme of the pipeline is the Leave One Patient Out cross-validation (LOPO). This scheme is a generalization of the classic leave one out cross-validation [10], with the difference that all data from the same patient are recursively left in the validation set, instead of a single data. This avoids biased high accuracy due to the similarity of the data coming from the same patient that may otherwise be found both in the training and validation set.

In Figure 1, we report the entirety of the dataset of Raman spectra, divided by the class of Alzheimer’s disease and the corresponding average with standard deviation.

In Figure 2 are shown eight Persistence images, two for each combination AD-noAD and RS-FT. It is interesting to note that in the PIs coming from the Raman spectra the pattern seems more chaotic between the two classes, while in the PIs coming from the Fourier transform there is a clearer division. More in

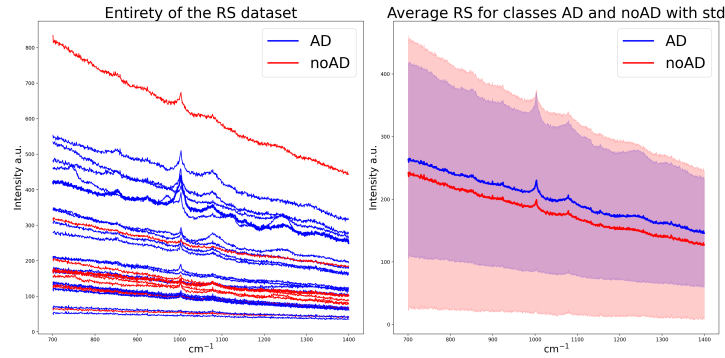


Fig. 1. (a) The entirety of the dataset of Raman spectra coloured by respective class. (b) The corresponding average with standard deviation.

detail, the lit pixels in the PIs of class noAD have a more elongated shape than those of the AD class. This corroborates the results achieved in Section 4. In Figure 3 are shown eight Persistence silhouettes in the same fashion as Figure 2. Again, there is a clearer division for the PSs coming from the Fourier transform. A peak at the tail end of the signal is present for the noAD class.

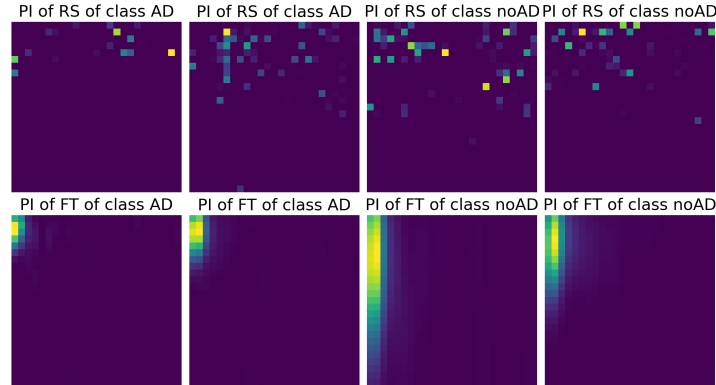


Fig. 2. First row: Two persistence images (PI) of the Raman spectra for class AD and two for class noAD. Second row: two PI of the Fourier transform for class AD and for class noAD. It appears that the PI obtained from Raman spectra are more chaotic between the two classes. In the PI obtained from the Fourier transform, a clearer division between AD and noAD is observed, with the latter having a more elongated dot.

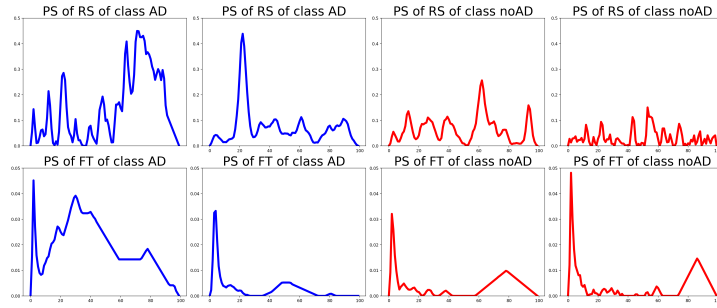


Fig. 3. First row: Two persistence silhouettes (PS) of the Raman spectra for class AD and for class noAD. Second row: two PS of the Fourier transform for class AD and for class noAD. Again, it seems that the PS coming from the RS are more chaotic, while in the PS coming from the Fourier transform there is a clearer division between AD and noAD, with the latter having a clear peak at the tail end of the signal.

4 Results

The results obtained from the pipeline on each of the transformations are shown in Table 1. The Fourier transform is the one that gets clearly the better results. We used as baseline accuracy the value of 0.733, due to the imbalance in classes (i.e. the accuracy achieved by the classifier assigning always the most frequent label to any sample).

Table 1. Accuracy result of the TML pipeline for different input.

Method	Accuracy	Vectorization and Classifier
H_0	0.833	PI and Ridge
Fourier transform	0.875	PS and SVC
Welch transform	0.763	PI and SVC
Autocorrelation	0.667	PI and Ridge

It is worth pointing out that even standard preprocessing applied to Raman spectra could lead to a classification accuracy below the baseline accuracy. This is probably due to the fact that in our dataset, the signal-to-noise ratio is quite low. On the other hand, the accuracy value ($> 83\%$) achieved by extracting H_0 features from raw spectra is in line with results of [14], while results achieved by extracting topological features after performing the Fourier transform are even better (87.5%).

5 Discussion

The results described above support strongly that RS and topological analysis together may provide an effective combination to confirm or disprove a clinical diagnosis of AD. Also, the training of the classification ML model trained on the topological features extracted from the Raman spectra acquired on CSF sample does not need the choice or set of any parameters; hence, the proposed methodology may evolve in automatic support to AD diagnosis, which could be easily embedded in a commercial platform of Raman spectroscopy. The above considerations are preliminary and require further confirmation from the statistical viewpoint. From this perspective, the next steps will include enlarging the dataset of CSF samples to validate the proposed method better and, possibly, to understand if topological machine learning could support the characterization of AD subtypes.

References

1. Adams, H., Emerson, T., Kirby, M., Neville, R., Peterson, C., Shipman, P., Chepushtanova, S., Hanson, E., Motta, F., Ziegelmeier, L.: Persistence images: A stable vector representation of persistent homology. *Journal of Machine Learning Research* **18** (2017)
2. Blennow, K., Zetterberg, H.: Biomarkers for alzheimer’s disease: current status and prospects for the future. *Journal of Internal Medicine* **284**(6), 643–663 (2018). <https://doi.org/https://doi.org/10.1111/joim.12816>
3. Breiman, L.: Random forests. *Machine learning* **45**(1), 5–32 (2001)
4. Bubenik, P., et al.: Statistical topological data analysis using persistence landscapes. *J. Mach. Learn. Res.* **16**(1), 77–102 (2015)
5. Chazal, F., Fasy, B.T., Lecci, F., Rinaldo, A., Wasserman, L.: Stochastic convergence of persistence landscapes and silhouettes. In: *Proceedings of the thirtieth annual symposium on Computational geometry*. pp. 474–483 (2014)
6. Conti, F., D’Acunto, M., Caudai, C., Colantonio, S., Gaeta, R., Moroni, D., Pascali, M.A.: Raman spectroscopy and topological machine learning for cancer grading. *Scientific reports* **13**(1), 7282 (May 2023). <https://doi.org/10.1038/s41598-023-34457-5>
7. Conti, F., Moroni, D., Pascali, M.A.: A topological machine learning pipeline for classification. *Mathematics* **10**(17) (2022). <https://doi.org/10.3390/math10173086>
8. Cortes, C., Vapnik, V.: Support-vector networks. *Machine learning* **20**(3), 273–297 (1995)
9. Eberhardt, K., Stiebing, C., Matthäus, C., Schmitt, M., Popp, J.: Advantages and limitations of raman spectroscopy for molecular diagnostics: an update. *Expert Review of Molecular Diagnostics* **15**(6), 773–787 (2015). <https://doi.org/10.1586/14737159.2015.1036744>, pMID: 25872466
10. Hastie, T., Tibshirani, R., Friedman, J.H., Friedman, J.H.: *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer (2009)
11. Hoerl, A.E., Kennard, R.W.: Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12**(1), 55–67 (1970)
12. Huang, C.C., Isidoro, C.: Raman spectrometric detection methods for early and non-invasive diagnosis of alzheimer’s disease. *Journal of Alzheimer’s Disease* **57**, 1145–1156 (2017). <https://doi.org/10.3233/JAD-161238>

13. Polykretis, P., Banchelli, M., D'Andrea, C., de Angelis, M., Matteini, P.: Raman spectroscopy techniques for the investigation and diagnosis of alzheimer's disease. *FBS* **14**(3), 22–null (2022). <https://doi.org/10.31083/j.fbs1403022>
14. Ryzhikova, E., Ralbovsky, N.M., Sikirzhytski, V., Kazakov, O., Halamkova, L., Quinn, J., Zimmerman, E.A., Lednev, I.K.: Raman spectroscopy and machine learning for biomedical applications: Alzheimer's disease diagnosis based on the analysis of cerebrospinal fluid. *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* **248**, 119188 (2021). <https://doi.org/https://doi.org/10.1016/j.saa.2020.119188>
15. Tashjian, R.S., Vinters, H.V., Yong, W.H.: *Biobanking of Cerebrospinal Fluid*, pp. 107–114. Springer New York, New York, NY (2019), https://doi.org/10.1007/978-1-4939-8935-5_11
16. Umeda, Y.: Time series classification via topological data analysis. *Information and Media Technologies* **12**, 228–239 (2017)
17. Vanderstichele, H., Bibl, M., Engelborghs, S., Le Bastard, N., Lewczuk, P., Molinuevo, J.L., Parnetti, L., Perret-Liaudet, A., Shaw, L.M., Teunissen, C., Wouters, D., Blennow, K.: Standardization of preanalytical aspects of cerebrospinal fluid biomarker testing for alzheimer's disease diagnosis: A consensus paper from the alzheimer's biomarkers standardization initiative. *Alzheimer's & Dementia* **8**(1), 65–73 (2012). <https://doi.org/https://doi.org/10.1016/j.jalz.2011.07.004>
18. Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman, K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, İ., Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P., SciPy 1.0 Contributors: *SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python*. *Nature Methods* **17**, 261–272 (2020). <https://doi.org/10.1038/s41592-019-0686-2>
19. Xu, Y., Pan, X., Li, H., Cao, Q., Xu, F., Zhang, J.: Accuracy of raman spectroscopy in the diagnosis of alzheimer's disease. *Frontiers in Psychiatry* **14** (2023). <https://doi.org/10.3389/fpsy.2023.1112615>