Al/ML for PFAS Contamination Source Attribution

Debbie Fagan, Jessie Yaros, Danish Hussain, Eva Brayfindley, Charlotte Roiger, and Tim J Johnson, PNNL

Introduction

Historic use of per- and polyfluoroalkyl substances (PFAS) and subsequent concerns about toxicity and carcinogenicity have led to efforts to limit exposure¹. Among the estimated 12,000 PFAS compounds that may be in use today, less than 100 can be identified analytically². Artificial Intelligence/Machine Learning (AI/ML) models capable of detecting both known and emerging contaminants using limited historical data are of keen interest. In particular, the ability to discriminate among PFAS sources, such as Aqueous Film Forming Foam (AFFF) and other commercial formulations (CF), is of interest to facilitate responsible environmental management. This study uses AI/ML methods on mass spectral data to achieve source attribution.

Data

Data for this study comes from the NIST PFAS database that uses the Database Infrastructure for Mass Spectrometry (DIMSpec) toolkit and contains LC-MS/MS spectra for 104 PFAS samples resulting in a total of 7,194 high-resolution MS/MS spectra (Figure 1). All MS/MS data are transformed into fixed-length numeric encodings by using intensity binning.

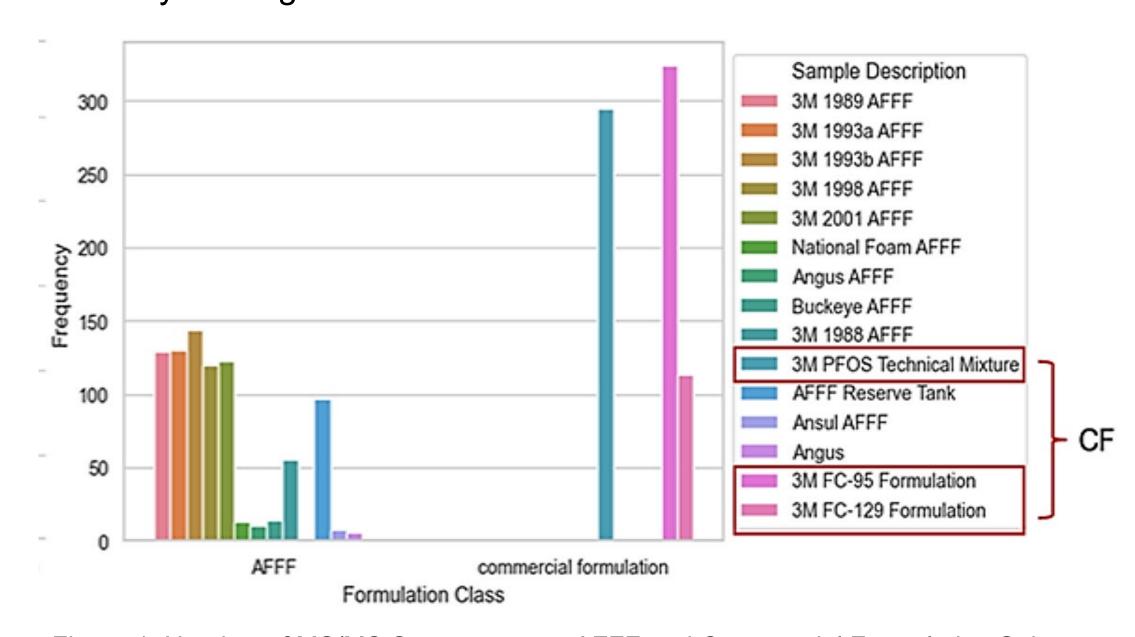


Figure 1. Number of MS/MS Spectra across AFFF and Commercial Formulation Subtype

For additional information, contact:



Debbie Fagan (509)-371-7784 | Deborah.fagan@pnnl.gov

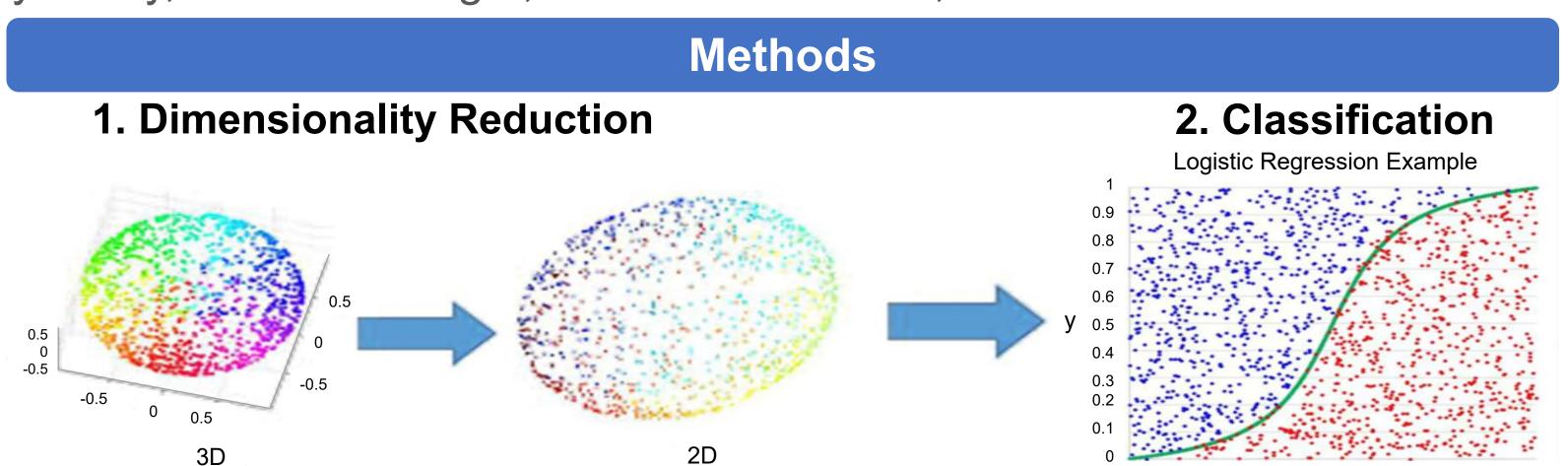


Figure 2. Schematic of Model Pipeline. A dimension reduction model is used, followed by a classification model

A series of three model pipelines were implemented with each consisting of a dimension reduction model followed by a classifier (Figure 2). Dimension reduction techniques include principal component analysis (PCA) and Uniform Manifold Approximation and Projection (UMAP). Three classifiers were employed: Linear Discriminant Analysis (LDA), Logistic Regression (LR), and Random Forests (RF). After testing multiple combinations of dimension reduction and classifiers, a combination of UMAP and RF models had the best performance for classifying AFFF and CF.

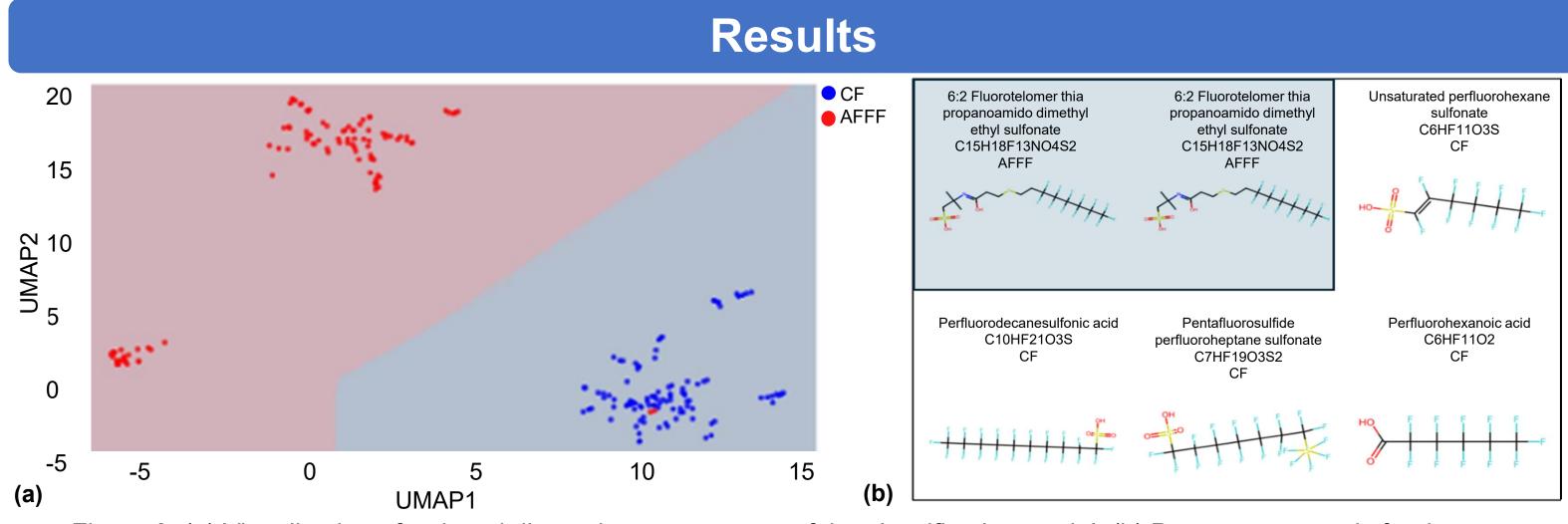


Figure 3. (a) Visualization of reduced dimension components of the classification model. (b) Parent compounds for the two misclassified AFFF MS2 fragments (boxed top left in grey) compared to other commercial formulations (CF).

After hyperparameter optimization, the model using UMAP and RF correctly distinguished AFFF from CF with 98% accuracy. Two AFFF samples were misclassified as CF, as shown in Figure 3(a). However, deeper investigation of sample compositions Figure 3(b) shows high similarity to other commercial formulations. Next steps will evaluate environmental samples with these models.

