SUPPLEMENTARY INFORMATION

Stochastic Generalization Models Learn to Comprehensively Detect Volatile Organic Compounds Associated with Foodborne Pathogens via Raman Spectroscopy

Bohong Zhang^{1,*}, Anand K Nambisan¹, Abhishek Prakash Hungund¹, Xavier Jones², Qingbo Yang^{2,*}, and Jie Huang^{1,*}

¹Department of Electrical and Computer Engineering, Missouri University of Science and Technology, Rolla, Missouri, 65409, USA

²Cooperative Research, College of Agriculture, Environmental and Human Sciences, Lincoln University of Missouri, Jefferson City, MO 65101, USA





Figure S1. Feature importance plots for (a) Rank-2 Random Forest and (b) Rank-1 XGBDT.

In **Table S1**, a variety of features are used to capture different characteristics of the signal. Features in this context include statistical metrics (e.g., mean, standard deviation, skewness), spectral peaks (e.g., intensity, position, width) derived from the wavenumber (Raman shift) domain, and frequency-domain characteristics (e.g., Fourier transforms, power spectral density). Specifically, the signal's frequency-domain features capture the power and logarithmic magnitude across different frequency scales, providing a deeper understanding of the signal's characteristics. These include the maximum, total, average, variance, and peak of the power and scaled logarithmic

magnitude. The skewness and kurtosis of the power and scaled logarithmic magnitude in the frequency domain are also calculated. These features might be correlated, particularly those related to power and magnitude. Additional complex features such as spectral roll-off, signal envelope features, Mel-Frequency Cepstrum Coefficients (MFCCs), and wavelet-transform-based features are also extracted. Spectral roll-off represents the frequency below which a given percentage of the total spectral energy is contained, useful for distinguishing between harmonic and non-harmonic content. Signal envelope features can include extracting the envelope of the signal using techniques like the Hilbert transform. MFCCs are coefficients that represent audio and are widely used in speech and audio processing. Wavelet transforms can also be used to extract features. Histogram bins are calculated for different feature sets, providing another view of the distribution of these features. The count of the values of those features within the range defined by the bins could also show correlation depending on the distribution of your data. In summary, these features provide a comprehensive representation of the signal, capturing its various characteristics in both the wavenumber and frequency domains. However, potential correlations between these features should be considered before the model pipeline training.

Table S1. Feature alias/definition table. This table contains the list of aliases used to represent each of the features used for training the classifiers. The aliases are assigned to each of the feature for ease of representation and to optimize the feature importance plotting. This table contains all the 96 extracted features for the model pipeline training. The correlation between each of these features is computed and analyzed in detail before getting to the model pipeline training.

Feature		Definition	Alias	
(a)	FREQ_SCALEDLOGMAG_TOP10_LOC_0	0th set of Top 10 peak locations of Scaled log -mag in frequency domain		
• •	PSD 23	23 rd Power spectral density		
	PSD_87	87th Power spectral density		
	RS_IQR	Interquartile range of the signal in wavenumber (Raman shift) domain, a measure of statistical dispersion	F_3	
	WAVELET_LEVEL_2_MEAN	Mean of the level-2 wavelet transforms	F_4	
	PSD_20	20th Power spectral density	F_5	
-	PSD_15	15th Power spectral density	F_6	
	PSD_70	70th Power spectral density	F 7	
	MFCC 9	9th Mel-Frequency Cepstrum Coefficients	F 8	
	PSD 64	64th Power spectral density	F 9	
	FREQ SCALEDLOGMAG BOT10 PKS 2	2 nd set of Bottom 10 peak magnitudes of Scaled log -mag in frequency domain	F 10	
	PSD 13	13 th Power spectral density	F 11	
	RS CREST FACTOR	Ratio of the peak value to the root mean square (RMS) value over a range wavenumbers	F 12	
	RS AUTOCORRELATION	Autocorrelation of the signal, measuring the linear relationship between lagged values of the signal	F 13	
	RS MEDIAN	Median value of the wavenumber (Raman shift) domain signal	F 14	
	FREQ HIST MAG 3	Magnitude of Histogram in 3 rd bin of frequency domain signal	F 15	
	SPECTRAL FLATNESS	Measure of the "flatness" or "tonality" of spectrum of the signal. A higher value indicates a more noise -like signal.	F 16	
	PSD 25	25th Power spectral density	F 17	
-	PSD 65	65 th Power spectral density	F 18	
	PSD 47	47 th Power spectral density	F 19	
	MECC 11	11th Mel-Erequency Censtrum Coefficients	F 20	
	WAVELET LEVEL 3 STD	Standard deviation of the level-3 wavelet transforms	F 21	
-	MFCC 17	17th Mel-Frequency Censtrum Coefficients	F 22	
	WN TOP5 PKS 3	3 rd set of wavenumbers associated with ton 5 neaks	F 23	
	PSD 32	37 nd Power spectral density	F 24	
	PSD 17	17 th Power spectral density	F 25	
	MECC 18	18th Mel-Frequency Censtrum Coefficients	F 26	
	WAVELET LEVEL 5 STD	Standard deviation of the level -5 wavelet transforms	F 27	
	EREO SCALEDLOGMAG BOTIO PKS 3	3rd set of Bottom 10 peak magnitudes of Scaled log -mag in frequency domain	F 28	
		Average scaled logarithmic manifulde in the frequency domain	F 29	
	PSD 76	74 th Power spectral density	F 30	
	PSD 7	7th Power spectral density	F 31	
	WAVELET LEVEL 1 MEAN	Mean of the level-1 wavelet transforms	F 32	
	FRED SCALEDLOGMAG TOP10 LOC 7	7th set of Ton 10 near 6 ratio of Scaled log - mar in frequency domain	F 33	
		60 th Dower spectral density	F 34	
-	HIST 2	2 rd high of Historian	F 35	
	WN TOP5 LOC 3	3rd set of wavenumber indires associated with ton 5 neaks	F 36	
	FREO HIST MAG 1	Magnitude of Histogram in 1% bin of frequency domain signal	F 37	
	WAVELET LEVEL 3 MEAN	Mean of the level-3 wavelet transforms	F 38	
	WN TOP5 LOC 1	1% set of wavenumber indices associated with ton 5 peaks	F 39	
	FRED SCALEDLOGMAG TOP10 LOC 3	3rd set of Transformations of Scaled log-map in frequency domain	F 40	
	MECC 5	5 th Mel-Frequency Censtrum Coefficients	F 41	
	PSD 66	6 th Power spectral depity	F 42	
	PSD 102	10 ²⁴ Power spectral density	F 43	
	PSD 34	34 th Power spectral density	F 44	
	PSD_28	28 th Power spectral density	F 45	
	PSD 72	7 ²⁰⁴ Power spectral density	F 46	
		r 2 Torrer spectral density	+0	

	MFCC_13	13th Mel-Frequency Cepstrum Coefficients	F_47	
h	PSD_96	96 ^h Power spectral density		
	WN_TOP5_PKS_2	2 nd set of wavenumbers associated with top 5 peaks		
	RS_ENTROPY	Entropy of the signal's probability distribution, representing the randomness or unpredictability.		
	PSD_101	101st Power spectral density		
	HIST_3	3rd bin of Histogram		
	FREQ_SCALEDLOGMAG_TOP10_PKS_4	4th set of Top 10 peak magnitudes of Scaled log -mag in frequency domain		
	FREQ_SCALEDLOGMAG_BOT10_LOC_5	5 th set of Bottom 10 peak locations of Scaled log -mag in frequency domain		
	PSD_81	81st Power spectral density		
	FREQ_SCALEDLOGMAG_TOP10_PKS_5	5th set of Top 10 peak magnitudes of Scaled log -mag in frequency domain		
	MFCC 10	10 th Mel-Frequency Cepstrum Coefficients		
	WAVELET LEVEL 5 MEAN	Mean of the level-5 wavelet transforms	F_58	
	PSD 14	14th Power spectral density		
	PSD 26	26 th Power spectral density	F 60	
	PSD 31	31st Power spectral density	F 61	
	HIST 1	1 st bin of Histogram	F 62	
	PSD 42	42 nd Power spectral density		
	WN TOP5 PKS 1	1 st set of wavenumbers associated with top 5 peaks	F 64	
	WN TOP5 LOC 4	4 th set of wavenumber indices associated with top 5 peaks	F 65	
	PSD 75	75 th Power spectral density		
	MFCC 1	1st Mel-Frequency Cepstrum Coefficients		
	FREQ VAR SCALEDLOGMAG	variance of the scaled logarithmic magnitude in the frequency domain	F 68	
	WAVELET LEVEL 4 MEAN	Mean of the level-4 wavelet transforms	F 69	
	PSD 97	97 th Power spectral density	F 70	
	MFCC 6	6 th Mel-Frequency Censtrum Coefficients		
	MFCC 12	12 th Mel-Frequency Censtrum Coefficients		
	WN TOP5 PKS 0	0 th set of wavenumbers associated with top 5 peaks		
	PSD 27	27 th Power spectral density		
	MFCC 15	15 th Mel-Erequency Centrum Coefficients		
	FREO SKEW POWER	skewness of the power in the frequency domain		
	PSD 100	100 th Power spectral density		
	PSD 68	68 th Power spectral density	F 78	
	PSD_61	61st Power spectral density	F 79	
	FREO HIST MAG 2	Magnitude of Histogram in 2 nd bin of frequency domain signal	F 80	
	PSD 99	99 th Power spectral density	F 81	
	WN TOP5 LOC 2	2 nd set of wavenumber indices associated with top 5 peaks	F 82	
	PSD 62	62 Set of waterhandle indices associated with top 5 peaks		
-	MECC 19	19th Mel-Frequency Censtrum Coefficients		
	MECC 8	8th Mel-Frequency Censtrum Coefficients	F 85	
	MECC 14	14th Mel-Frequency Censtrum Coefficients	F 86	
	MFCC 16	16 th Mel-Frequency Cepstrum Coefficients	F 87	
-	SPECTRAL ROLLOFF	frequency below which a given percentage (e.g. 85%) of the total spectral energy is contained	F 88	
	WN TOP5 LOC 0	Oth set of wavenumber indices associated with ton 5 peaks	F 89	
	MFCC 7	7th Mel-Frequency Censtrum Coefficients		
	MFCC 2	2014 Mel-Frequency Cepstain Coefficients		
	MECC 0	Othel-Frequency Censtrum Coefficients	F 92	
	MFCC 3	3 rd Mel-Frequency Censtrum Coefficients	F 93	

For all the features given in the **Figure S1 and S2**, the correlations are calculated w.r.t. one another and represented in a matrix form. The features showing lower correlations are represented by white or lighter shade of red and the features showing higher correlations are represented as darker shade of red. The purpose of this calculation is to determine the features having a high level of correlation. Those features can affect the model pipeline performance due to the multi-collinearity causing inaccuracies. Such highly correlated features are dropped during the preprocessing before the pipeline training step. The threshold to drop features based on correlation analysis is set to 0.95.



Figure S2. Feature correlation matrix.

- 0.8



Figure S3. Regression to predict the dilution levels of all the VOCs. VOC dilution level predictions by (a) Random Forest rank-2 regression and (b) XGBDT rank-1 decision tree regression.



(a) Rank-1 MLP Classifier



(b) Rank-3 Random Forest Classifier



(c) Rank-2 XGBDT Classifier

Figure S4. Feature importance plots for 3 prominent pipeline classifiers. FI plot for (a) rank-1 MLP classifier, (b) rank-3 random forest classifier, and (c) rank-2 XGBDT classifier. The plots are tabulated according to the pure VOCs' class labels S1 to S8. When these VOCs are mixed to obtain a mixture VOC, the feature importance (or mean of SHAP values) of that mixture will be a derivative of the mean SHAP value of the pure VOCs used to obtain the mixture. The feature aliases in the x-axis of the bar plots are defined in the Table. S1 in the Supplementary Information. These plots give detailed insight into the impact of certain features on the model pipeline performance and are plotted right to left with the most important feature (i.e., the feature with largest mean SHAP value) starting at the right. Higher mean SHAP values suggest a greater impact on the model's output. Features with larger bars are more influential in the model's predictions.



Figure S5. Raman spectra of eight pure VOCs and ACN.

Dilution Level (X)	Concentration (unit/mL)	Concentration (%)	Concentration (PPM)
OX (Pure)	1	100.00%	1,000,000
5X	0.2	20.00%	200,000
10X	0.1	10.00%	100,000
20X	0.05	5.00%	50,000
40X	0.025	2.50%	25,000
60X	0.0167	1.67%	16,667
80X	0.0125	1.25%	12,500
100X	0.01	1.00%	10,000
160X	0.0063	0.63%	6,250
200X	0.005	0.50%	5,000
240X	0.0042	0.42%	4,167
320X	0.0031	0.31%	3,125
400X	0.0025	0.25%	2,500

Table S2. Concentration Levels of VOC Mixtures at Various Dilution Factors Expressed in Units per Milliliter (unit/mL), Percentage (%), and Parts Per Million (PPM).