

Bittersweet

Manveet Singh², Harsh Hingorani², Karan Yadav², Mansi Goel^{1,3,4}, and
Ganesh Bagler^{1,3,4}

^a*Department of Computational Biology, Indraprastha Institute of Information Technology
Delhi (IIIT-Delhi), New Delhi, 110020, India*

^b*Department of Computer Science, Indraprastha Institute of Information Technology
Delhi (IIIT-Delhi), New Delhi, 110020, India*

^c*Infosys Centre for Artificial Intelligence, Indraprastha Institute of Information
Technology Delhi (IIIT-Delhi), New Delhi, 110020, India*

^d*Center of Excellence in Healthcare, Indraprastha Institute of Information Technology
Delhi (IIIT-Delhi), New Delhi, 110020, India*

Abstract

Abstract text.

Keywords:

1. Introduction

Humans use their fundamental sense of taste to select food while avoiding toxic substances through an evolutionary process. Among the five basic tastes, sweetness and bitterness represent contrasting cues: sweetness indicates energy-rich carbohydrates, thus attracting individuals, while bitterness warns about dangerous substances, which lead to aversion [1]. G-protein-coupled receptors—T1R2/T1R3 and approximately 25 hTAS2Rs act as mediators for the detection of sweet and bitter substances in the body [2]. These taste receptors extend their presence beyond the mouth to detect chemical signals in different tissues to regulate vital bodily functions, including appetite control and metabolism, and glucose management [3].

Over the last few years, machine learning (ML) has proven to be an influential tool in the field of computational gastronomy that can predict and control taste profiles according to molecular properties. Classical experimental methods are hampered by ethical issues, high expenses, and inefficiency in processing big compound libraries [4]. ML provides scalable solutions that automate the identification of new sweeteners and bitter-maskers, especially

18 in the food, pharmaceutical, and health industries [5, 6]. More recent studies
19 have shown the effectiveness of ML models such as random forests, neural
20 networks, and deep learning models in the prediction of flavor profiles from
21 chemical and structural information [7, 8].

22 Most of the previous research has focussed upon binary classification,
23 classifying between sweet/non-sweet or bitter/non-bitter molecules [9, 10].
24 But there are not many that address the complete bitter-sweet continuum or
25 investigate the common and unique chemical characteristics of both flavors.
26 Additionally, most current models depend on fewer molecular descriptors or
27 do not consider newer methods such as convolutional neural networks (CNNs)
28 and transformer models, which have demonstrated significant potential in
29 adjacent fields such as flavor chemistry and smell [11, 12].

30 We revisited this problem to improve upon the existing research by com-
31 bining a broader set of molecular descriptors which allowed more precise clas-
32 sification across the bitter-sweet-tasteless spectrum, and offering considerable
33 insights into molecular taste properties. To solve the existing limitations and
34 challenges, we developed our Bitter-Sweet Taste Prediction pipeline using a
35 two-way approach.

36 We initially trained machine learning models—Random Forest (RF) and
37 AdaBoost (AB)—on Boruta-selected molecular descriptors to classify com-
38 pounds as bitter/non-bitter and sweet/non-sweet. We next trained neural
39 network models—Neural Networks (NN) and Convolutional Neural Networks
40 (CNN)—on multi-class classification as bitter, sweet, or tasteless using com-
41 bined taste datasets. The NN model had four dense layers, while the CNN
42 model had two Conv1D layers, with dropout and early stopping techniques
43 to prevent overfitting.

44 2. Materials and Methods

45 2.1. Data Compilation and Curation

46 In this study, we compiled and curated data from a previous scientific
47 study, BitterSweet, which carefully compiled a dataset from a variety of
48 resources and excluded molecules with unverified or conflicting taste infor-
49 mation. The following resources were used to compile data: (a) Biochemical
50 Targets of Plant Bioactive Compounds by Gideon Polya; (b) BitterDB; (c)
51 Fenaroli’s Handbook of Flavor Ingredients (5th Edition); (d) Rodgers et al.;
52 (e) Rojas et al.; (f) SuperSweet; (g) TOXNET; (h) The Good Scents Com-
53 pany Database (www.thegoodscentscompany.com); and (i) BitterPredict.

54 The data set consists of 918 bitter and 1510 non-bitter molecules and 1247
 55 sweet and 1119 non-sweet molecules. For the creation of machine learning
 56 models, the data set was split into training and testing. The data set consists
 57 of the name and structure of molecules stored in string format using SMILES
 58 (Simplified Molecular Input Line-Entry System) along with bitter/non-bitter
 59 and sweet/non-sweet classification. To visualize the high dimensional feature
 60 space training dataset, 2D t-SNE plot was generated using all set of molecular
 61 descriptors (Fig 2). The complete schema describing the process of data
 62 collection, feature selection, and model building is depicted in Fig 1

63 2.2. Feature Generation

64 Molecule descriptors define the physical, chemical, and structural prop-
 65 erties of a molecule. Online Chemical Modeling Environment (see- <https://ochem.eu/login/show.do>)[13] was used to calculate 15046 molecular de-
 66 scriptors, namely Dragon 2D/3D (5270), ECFP (1024), JPlogP (115), MOLD2
 67 (777), MORDRED (1613), PaDEL2 (2756), and RDKit (3941).

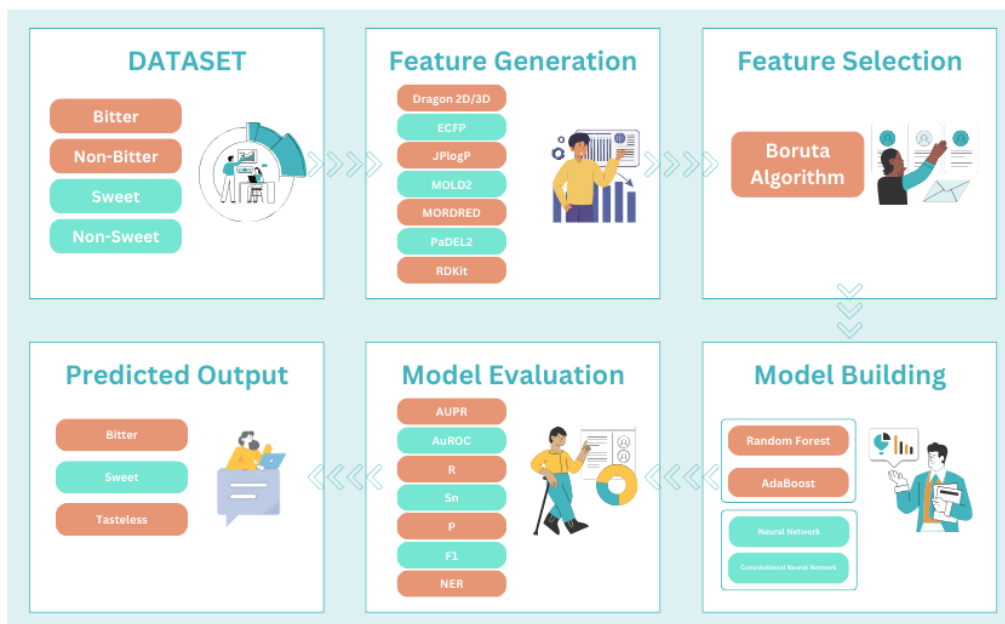


Figure 1: Schema for building bitter/non-bitter and sweet/non-sweet prediction models

69 2.3. Feature Selection

70 Since lot of descriptors were common among different classes of molecular
71 descriptors, we removed duplicates and imputed null/missing values with
72 mean of the concerned feature.

73 Standard-Scaler was used to normalize each descriptor. The outliers were
74 removed using a Quantile transformer. Features having low variance and high
75 correlation were removed using Variance-Threshold method and Correlation-
76 Based method to avoid redundancy and obtain relevant descriptors.

77 Further Boruta algorithm was applied to remove irrelevant features. A
78 significant proportion of features (65%) in the remaining descriptor were
79 identified as irrelevant by the Boruta algorithm.

80 2.4. Model Implementation

81 In order to classify a molecule into bitter, sweet and tasteless, two different
82 approaches were implemented. Firstly, Random Forest (RF) and Adaboost
83 (AB) machine learning models were implemented to classify a molecule as
84 bitter/non-bitter and sweet/non-sweet, using molecular descriptors chosen
85 from the Boruta algorithm.

86 Secondly, Neural Network (NN) and Convolutional Neural Network (CNN)
87 deep learning models were implemented to classify a molecule as bitter, sweet
88 or tasteless. Bitter and sweet data was integrated into single dataset. NN
89 with four hidden layers of 512, 256, 128, and 64 neurons was used. CNN
90 with two Conv1D layers (64 & 128 filters, kernel size 3) was used. For regu-
91 larization and to avoid overfitting, dropout method and early stopping were
92 implemented, respectively. For measuring model performance, threshold-
93 independent measures were employed to prevent classification inconsistencies:
94 Area Under the Precision-Recall Curve (AUPR), Area Under the Receiver
95 Operating Characteristic Curve (AUROC) F1-score, Specificity, Sensitivity.
96 Further to find features which contributed most to the prediction of molecule
97 into bitter/non-sweet and sweet/non-sweet SHAP (SHapley Additive exPla-
98 nations) Analysis was used.

99 *Random forest (RF)*. Random Forest algorithm is a type of ensemble
100 learning method that utilizes bagged decision trees. They are quite versatile
101 and can be used for both classification as well as regression. RF works by
102 building a number of decision trees (usually greater than 100) at training
103 time, each utilizing a subset of features and data points. At the time of pre-
104 diction, the predictions made by its constituent decision trees are aggregated



Figure 2: 2D t-SNE scatterplot generated using all set of molecular descriptors for bitter/non-bitter and sweet/non-sweet.

105 *Adaboost (AB)*. AdaBoost is an ensemble method that operates by itera-
 106 tively training weak classifiers in sequence, each of which puts more empha-

sis on the instances previously misclassified by the earlier ones. It assigns weights to the training examples, raising them for the misclassified samples so that the next classifier puts more emphasis on them. It then pools all weak classifiers by a weighted majority vote (in classification) or weighted sum (in regression).

Neural Network (NN). A Neural Network is a machine learning model inspired by the human brain. It is made up of layers of nodes (neurons) that are interconnected and feed data. NNs learn patterns by adjusting the weights with the assistance of backpropagation. NNs are applied in classification, regression, and pattern recognition tasks.

Convolutional Neural Network (CNN). A Convolutional Neural Network (CNN) is a form of deep learning architectures that are particularly designed to process and understand grid-structured data, e.g., images. CNNs can automatically and flexibly learn spatial hierarchies of features and are able to do so through the utilization of convolutional layers that perform filtering on the input data, thereby extracting meaningful patterns such as edges, textures, or objects.

SHAP (SHapley Additive exPlanations) is one of the most widely used explainable AI (XAI) methods, which provides a contribution value (a SHAP value) to every feature of a machine learning model to interpret the prediction for an instance. The concept relies on cooperative game theory's Shapley values, in which the prediction is considered a "payout" shared equally among all the features according to their individual contribution.

3. Results

Model	AUPR	AUROC	F1	NER	Sn	Sp
RF Model	0.9175	0.9686	0.8631	0.9006	0.9002	0.9394
CNN Model	0.9162	0.9357	0.8515	0.9048	0.8571	0.8821
NN Model	0.9014	0.9326	0.8616	0.9089	0.8634	0.9034
AB Model	0.5841	0.7630	0.6249	0.7453	0.6883	0.6667

Table 1: Comparison of performance on the test sets of Random Forest (RF) model, AdaBoost (AB) Model, Convolutional Neural Network (CNN) and Neural Network (NN) model

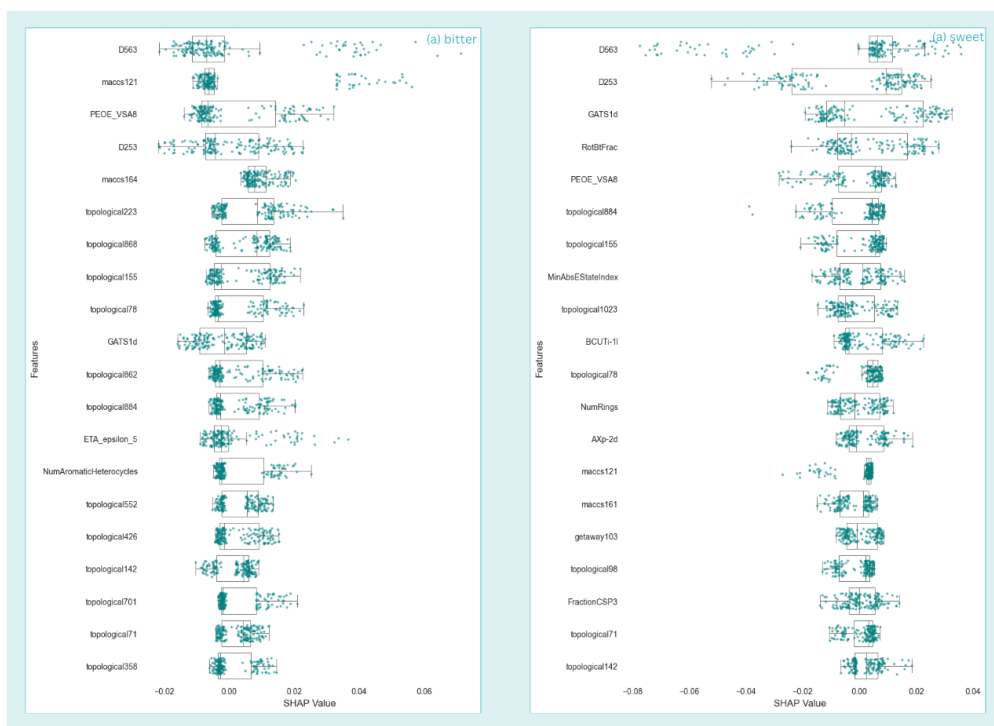


Figure 3: Box plot of importance scores of molecular descriptors for (a) bitter/non-bitter prediction and (b)sweet/non-sweet prediction.

4. Discussion

5. Comparison with existing studies

6. Conclusion

References

- [1] A. B. Morrison, Sugar substitutes, Canadian Medical Association Journal 120 (6) (1979) 633–637.
- [2] A. Di Pizio, M. Y. Niv, Computational studies of smell and taste receptors, Israel Journal of Chemistry 54 (8–9) (2014) 1205–1218. doi:10.1002/ijch.201400027.
- [3] M. S. Bahia, I. Nissim, M. Y. Niv, Bitterness prediction in-silico: A step towards better drugs, International Journal of Pharmaceutics 536 (2) (2018) 526–529. doi:10.1016/j.ijpharm.2017.03.076.

- 143 [4] A. Levit, S. Nowak, M. Peters, A. Wiener, W. Meyerhof, M. Behrens,
144 M. Y. Niv, The bitter pill: Clinical drugs that activate the human bitter
145 taste receptor tas2r14, *The FASEB Journal* 28 (3) (2014) 1181–1197.
146 [doi:10.1096/fj.13-242594](https://doi.org/10.1096/fj.13-242594).
- 147 [5] A. Sharma, S. Amarnath, M. Thulasimani, S. Ramaswamy, Artificial
148 sweeteners as a sugar substitute: Are they really safe?, *Indian Journal*
149 *of Pharmacology* 48 (3) (2016) 237. [doi:10.4103/0253-7613.182888](https://doi.org/10.4103/0253-7613.182888).
- 150 [6] H. Ji, D. Pu, W. Yan, Q. Zhang, M. Zuo, Y. Zhang, Recent advances and
151 application of machine learning in food flavor prediction and regulation,
152 *Trends in Food Science & Technology* 138 (2023) 738–751. [doi:10.](https://doi.org/10.1016/j.tifs.2023.07.012)
153 [1016/j.tifs.2023.07.012](https://doi.org/10.1016/j.tifs.2023.07.012).
- 154 [7] M. Schreurs, S. Piampongsant, M. Roncoroni, L. Cool, B. Herrera-
155 Malaver, C. Vanderaa, F. A. Theßeling, L. Kreft, A. Botzki, P. Malcorps,
156 L. Daenen, T. Wenseleers, K. J. Verstrepen, Predicting and improving
157 complex beer flavor through machine learning, *Nature Communications*
158 15 (1) (2024) 2368. [doi:10.1038/s41467-024-46346-0](https://doi.org/10.1038/s41467-024-46346-0).
- 159 [8] T. Naravane, I. Tagkopoulos, [Machine learning models to predict mi-](https://doi.org/10.1016/j.crfs.2023.100498)
160 [cronutrient profile in food after processing](https://doi.org/10.1016/j.crfs.2023.100498), *Current Research in Food*
161 *Science* 6 (2023) 100498. [doi:10.1016/j.crfs.2023.100498](https://doi.org/10.1016/j.crfs.2023.100498).
162 URL <https://doi.org/10.1016/j.crfs.2023.100498>
- 163 [9] W. Bo, D. Qin, X. Zheng, Y. Wang, B. Ding, Y. Li, G. Liang, Predic-
164 tion of bitterant and sweetener using structure-taste relationship models
165 based on an artificial neural network, *Food Research International* 153
166 (2022) 110974. [doi:10.1016/j.foodres.2022.110974](https://doi.org/10.1016/j.foodres.2022.110974).
- 167 [10] W. Huang, Q. Shen, X. Su, M. Ji, X. Liu, Y. Chen, S. Lu, H. Zhuang,
168 J. Zhang, Bitterx: A tool for understanding bitter taste in humans,
169 *Scientific Reports* 6 (2016) 23450. [doi:10.1038/srep23450](https://doi.org/10.1038/srep23450).
- 170 [11] L. Androutsos, L. Pallante, A. Bompotas, F. Stojceski, G. Grasso,
171 D. Piga, G. Di Benedetto, C. Alexakos, A. Kalogeras, K. Theofilatos,
172 M. A. Deriu, S. Mavroudi, Predicting multiple taste sensations with
173 a multiobjective machine learning method, *npj Science of Food* 8 (1)
174 (2024) 47. [doi:10.1038/s41538-024-00287-6](https://doi.org/10.1038/s41538-024-00287-6).

- 175 [12] L. Pallante, A. Korfiati, L. Androutsos, F. Stojceski, A. Bompotas,
176 I. Giannikos, C. Raftopoulos, M. Malavolta, G. Grasso, S. Mavroudi,
177 A. Kalogeras, V. Martos, D. Amoroso, D. Piga, K. Theofilatos, M. A.
178 Deriu, Toward a general and interpretable umami taste predictor using
179 a multi-objective machine learning approach, Scientific Reports 12 (1)
180 (2022) 21735. doi:10.1038/s41598-022-25935-3.
- 181 [13] Online chemical modeling environment, [https://ochem.eu/login/
182 show.do](https://ochem.eu/login/show.do), accessed April 11, 2025.