# Classification-based Machine Learning Approaches to Predict the Taste of Molecules: A Review

Cristian Rojas[a,*], Davide Ballabio[b], Viviana Consonni[b], Diego Suárez-Estrella[a], Roberto Todeschini[b]

[a] Grupo de Investigación en Quimiometría y QSAR, Facultad de Ciencia y Tecnología, Universidad del Azuay, Av. 24 de Mayo 7-77 y Hernán Malo, Cuenca 010107, Ecuador

[b] Milano Chemometrics and QSAR Research Group. Department of Earth and Environmental Sciences, University of Milano-Bicocca, P.za della Scienza 1-20126, Milano, Italy

[*]Corresponding author. E-mail: crojasvilla@gmail.com

## Abstract

The capacity to discriminate safe from dangerous compounds has played an important role in the evolution of species, including human beings. Highly evolved senses such as taste receptors allow humans to navigate and survive in the environment through information that arrives to the brain through electrical pulses. Specifically, taste receptors provide multiple bits of information about the substances that are introduced orally. These substances could be pleasant or not according to the taste responses that they trigger. Tastes have been classified into basic (sweet, bitter, umami, sour and salty) or non-basic (astringent, chilling, cooling, heating, pungent), while some compounds are considered as multitastes, taste modifiers or tasteless. Classification-based machine learning approaches are useful tools to develop predictive mathematical relationships in such a way as to predict the taste class of new molecules based on their chemical structure. This work reviews the history of multicriteria quantitative structure-taste relationship modelling, starting from the first ligand-based (LB) classifiers proposed in 1980 by Lemont B. Kier and concluding with the most recent studies published in 2022.

## Keywords

Taste chemistry; Machine learning; Taste classification; QSAR models; Foodinformatics

## 1. Background

## 1.1. Taste chemistry

### 1.1.1. Introduction

Considering the incredible variability of environmental conditions on the planet, the availability of specific foods has played a key role in the adaptive evolution and conservation of species. Indeed, the availability of specific types of nutrition may be one of the most important variables in the evolution of species. Taste and olfaction are the two senses that allow the discrimination of chemical substances (Schieberle & Hofmann, 2016). Dangerous tastes have been empirically correlated with bitterness; however, some of the most ancient medicines include bitter substances (Bayer *et al.*, 2021). Recently, scientific approaches have been replacing empiric ways to understand and to assess the safety of food products. While these scientific approaches have been shown to be useful in the analysis and categorizing of tastes, some foods that are considered safe may illicit, intolerances and allergies in a few people. For instance, human intolerances to gluten and lactose are well known, along with the life threatening anaphylactic allergic response to seafood and peanuts. These relatively rare reactions to food are related to specific digestive enzymes concentrations, and to the human immune system, respectively.

Beyond safety considerations, each person has specific preferences for tastes that can change over one's lifetime. This variation in personal taste preferences could be related to biochemical, as well as environmental, psychological and cultural factors. This diversity of factors makes it difficult to describe the taste mechanisms, e.g. psychologists might consider taste preferences to be related mainly to psychological stimulus, while chemists might think that tastes are perceived primarily through the consequence of chemical reactions that occur in tissues (Behrens & Ziegler, 2020), contributing to kind and intensity of sensations. In this framework, the identification of the healthiest, safest and most preferred foods is of fundamental importance to the food and pharmaceutical industries. Research related to the mechanisms underlying the human perception of tastes is increasing in the last few years (Damodaran & Parkin, 2017). Importantly, current research into tastes and sensory perceptions are being studied from different, but related and interconnected directions, such as chemical, biochemical, anatomical, physiological, and psychological standpoints.

It is common to identify "tastes" or "flavors" as the combination of taste, olfactory, tactile and thermal sensations (Di Lorenzo *et al.*, 2009), while sensomics is the mapping of the combinatorial code of aroma and taste by active key molecules. These molecules are sensed by human chemosensory receptors (Schieberle & Hofmann, 2016). The extraordinary developments in foodinformatics (computational food chemistry) and bioinformatics (computational biochemistry) are providing new

tools to assess and to explain the receptor/ligand binding affinity and how the structures of the receptors interact with the chemical structures of the compounds and how to achieve a particular taste of interest (Martinez-Mayorga & Medina-Franco, 2014; Rojas *et al.*, 2016a). At the beginning of sensory research, the greatest efforts were focused on the chemical structure of compounds, their characteristics and the cultural particularities of populations. It was considered that chemical analysis of taste molecules in raw ingredients and in end-products for human consumption could play an important role for the assurance of food quality and desirability preventing defects in products (Ley *et al.*, 2012). However, this approach was not sufficient to explain all the taste phenomena. Later, the importance of the complex anatomy and physiology of taste receptors and how they interact with specific tastant molecules were recognized as key factors to better understand and model the phenomenon of taste.

A molecular tastant is considered to be a water-soluble chemical able to produce taste sensations by activating taste receptor cells (TRCs) and thus activate taste-related pathways at within the nervous system (Di Lorenzo *et al.*, 2009; Rojas *et al.*, 2022). Tastants are elicited not only in water, but also in organic and inorganic acids and amino acids, all of which are able to facilitate the interactions of tastants with receptors (Chaudhari *et al.*, 2009). Chemosensory receptors located in the taste buds of the tongue are fundamental to the regulation of taste sensation. Other mechanisms to recognize molecular tastants are, for example, the opening of ion channels or through secondary messenger channels associated with nucleotides or phosphorylated inositol (Damodaran & Parkin, 2017; Morini *et al.*, 2011; Wong, 2018).

Taste measurement is preferably performed by an experienced panel of assessors. Panelists are trained with standard solutions of the basic tastes by means of the sip and spit methodology (Kelly *et al.*, 2005; Spillane *et al.*, 2006). The concentrations of standard solutions should be prepared at a minimum of their recognition threshold to ensure taste detection (Deng *et al.*, 2021; Liu *et al.*, 2020; Shiyan *et al.*, 2021; Spillane *et al.*, 2006; Yu *et al.*, 2018). The pH of the standard solutions also influences taste perception. Then, a solution of an unknow analyte (generally at concentration of 0.01 M) is provided to members of the panel who are asked to identify the basic taste and aftertaste. The taste potency of the unknown analyte can be estimated by the amount that the solution should be diluted to be equal to the standard. The evolution of technology led to the development of some analytical procedures based on sensors for the sensory evaluation of foods; for instance, electronic noses and tongues, in which their operation is based on the measurement of potential differences that are related to the tastes and aromas that humans can sense (Deng *et al.*, 2021; Liang *et al.*, 2022a; Suárez-Estrella *et al.*, 2021; Xiu *et al.*, 2022).

### 1.1.2. Basic tastes

Currently, five basic tastes have been identified: sweet, bitter, umami, sour and salty, which are referred to as basic taste modalities, taste qualities or receptor-mediated tastes (Chandrashekar *et al.*, 2006; Damodaran & Parkin, 2017; Di Lorenzo *et al.*, 2009; Morini *et al.*, 2011; Wong, 2018). Among the basic tastes, sweetness is probably the most important one, since sweeteners evoke a high caloric intake and a pleasant sensation in many foods and medicines (Chandrashekar *et al.*, 2006; Damodaran & Parkin, 2017). Most sweet foods contain mono- and disaccharides (Di Lorenzo *et al.*, 2009), which are responsible for their sweetness and quick sources of energy for the body. On the other hand, several non-caloric substances capable of providing a sensation of sweetness to food are currently known and are used in the industry. Those substances may have both natural or artificial origins (Chattopadhyay *et al.*, 2014).

Sweetness perception is related to the presence of a glycophore unit in the sweetener's scaffold. It forms the tripartite model (AH, B and γ units), which interacts with the sweetness receptor along a multipoint attachment (MPA) construct. The sweet taste chemoreceptor is a G-protein coupled receptor (GPCR) of class C composed of the T1R2 (Type 1 Receptor 2) and T1R3 (Type 1 Receptor 3) subunits, which are composed of three structural domains (Chandrashekar *et al.*, 2006; Morini *et al.*, 2011; Wong, 2018). The presence of the AH-B sites in a tastant molecule is a necessary, but not a sufficient condition alone to elicit sweetness; for example, the sweetness taste can be viewed as a function of the size, shape and functionality of the compounds (Spillane & Sheahan, 1989). In other words, a large molecule must be able to fit specifically into the receptor site to generate sweetness. A small molecule with the AH-B site might be unable to match the construct of the receptor site, and the sweet stimulus may not be produced.

Sucrose has a clean (no aftertastes) sweet sensation (even at high concentrations), and consequently it is frequently used as the standard to quantify the relative sweetness (RS) or sweetness potency (Sw) of sweet-tasting molecules (Liu *et al.*, 2020; Rojas *et al.*, 2022; Shiyan *et al.*, 2021; Yu *et al.*, 2018). Sweet potency is defined as the concentration ratio between a sucrose solution standard labeled as 1 (or 100%), and the solution of a sweetener exhibiting the same intensity (iso-sweet concentration) (Rojas *et al.*, 2016a; Rojas *et al.*, 2016b). Sweeteners could be classified as natural (nutritive or carbohydrate) and artificial (non-nutritive or non-carbohydrate) (Ley *et al.*, 2012; Wong, 2018; Yang *et al.*, 2022). On the other hand, certain amino acids and proteins are detected as sweet compounds and some salts taste sweet at low concentrations, including NaCl, KCl, NaOH, KOH, salts of beryllium and lead acetate and carbonate (Di Lorenzo *et al.*, 2009).

Bitterness has been defined as an unpleasant taste. The unpleasant sensation is related to a rejection of some foods, many of which are toxic compounds for humans (Chandrashekar *et al.*, 2006; Di

4

129   Lorenzo *et al.*, 2009). Thus, bitter perception might be related to an evolved "alert" system to prevent
130   the intake of high concentration of toxic compounds through food or drink, avoiding their undesirable
131   and potential lethal effects (Ley *et al.*, 2012). On the other hand, not all bitter compounds are toxic
132   and not all toxic compounds are bitter. In fact, some of them have proven beneficial effects for human
133   health, for instance, polyphenols, glucosinolates and terpenes (Bayer *et al.*, 2021). Moreover, some
134   bitter tastes may be perceived as pleasant (Dagan-Wiener *et al.*, 2019) as well as associated food
135   products, such as coffee, beer, olives, and unsweetened chocolate. Plants that are perceived as slightly
136   bitter are commonly used for food, while plants perceived as highly bitter are more commonly used
137   in medicines. Plants perceived as having intermediate bitterness might be used for both alimentation
138   and/or medical purposes (Pieroni *et al.*, 2007; Pieroni *et al.*, 2002).

139   Bitter molecules generally require the presence of a polar (electrophilic or nucleophilic) group and a
140   hydrophobic group to interact with the bitter receptor. Bitter taste stimuli are associated with 25
141   receptors (TAS2Rs), which are G protein-coupled (Adler *et al.*, 2000; Chandrashekar *et al.*, 2006;
142   Matsunami *et al.*, 2000). Most of them are located in the same taste receptor cells (TRCs)
143   (Chandrashekar *et al.*, 2006; Damodaran & Parkin, 2017; Di Lorenzo *et al.*, 2009; Wong, 2018).
144   TAS2Rs have not only been identified in the mouth cavity, but also in gastrointestinal, respiratory,
145   reproductive and urinary tract tissues (Bayer *et al.*, 2021). The physiological function of TAS2Rs
146   outside the oral cavity have not been identified. Bitter receptors can be specific for one or a few
147   compounds, while others are able to react to a large number of bitter substances (Di Pizio & Niv,
148   2015). Some bitter compounds are agonists for some TAS2R subtypes, but antagonists for others
149   (Brockhoff *et al.*, 2011). Bitterness is a common taste reaction to alkaloids and heavy metal salts.
150   Quinine sulfate is the standard used for comparisons among the bitterness of compounds (Dagan-
151   Wiener *et al.*, 2017; Damodaran & Parkin, 2017). Quinine sulfate (Liu *et al.*, 2020; Rojas *et al.*, 2022)
152   and L-isoleucine (Shiyan *et al.*, 2021; Yu *et al.*, 2018) are the most frequently used standards for
153   bitterness identification. Quinine is an alkaloid used in the food industry as a component of some soft
154   drinks to infuse them with bitter taste, for example, tonic water. Substances used as sweeteners, such
155   as sodium saccharine and acesulfame K can become bitter at high concentration and also produce a
156   bitter aftertaste (Di Lorenzo *et al.*, 2009).

157   Umami is the most recently recognized basic taste. Umami is a Japanese word that means
158   deliciousness. This taste is associated with L-amino acids (such as monosodium glutamate MSG),
159   that are umami enhancers (potentiators) (Baines & Brown, 2016; Damodaran & Parkin, 2017; Suess
160   *et al.*, 2015; Wong, 2018). For instance, MSG exhibits a synergistic effect (enhancement) with the
161   guanosine 5'-monophosphate or inosine 5'-monophosphate nucleotides, although these compounds
162   also show a weak intrinsic umami taste on their own (Ley *et al.*, 2012; Wong, 2018). Also, L-aspartate

163 produces an umami sensation. Umami taste is detected in meats, cheeses, some mushrooms along
164 with fish, kelp and tomatoes. The umami taste stimuli of peptides and their molecular interactions is
165 associated with G-protein coupled receptors (GPCRs) comprised of the subunits T1R1 (Type 1
166 Receptor 1) and T1R3 (Type 1 Receptor 3) (Liang *et al.*, 2022a; Liang *et al.*, 2022b; Morini *et al.*,
167 2011). An umami nucleotide binds with the corresponding receptor at three points: two of them are
168 electrophilic (A and B) that interact with the two phosphoryl oxygens and the C6 oxygen,
169 respectively, while site X interacts with the substituent at C2, particularly when the substituent is
170 delocalized (Wong, 2018). The standard used to quantify umami intensity is MSG (Baines & Brown,
171 2016; Liu *et al.*, 2020; Rojas *et al.*, 2022; Shiyan *et al.*, 2021; Yu *et al.*, 2018).

172 Sour taste is associated with the presence of organic and inorganic acids in food. Acidity in raw food
173 tends to change with time; for example, acidity in soft fruit decrease as the fruit becomes ripe. A sour
174 taste is associated with unripe soft fruit. Sourness increases also after fermentation processes applied
175 for the production of foods, such as yogurt, wine, vinegar and bread. Initially, sourness perception
176 was related to the capacity of substances to release hydrogen ions in water. However, hydrogen ion
177 release is not the mechanism that produces sourness for organic and diluted inorganic acids (Breslin
178 & Huang, 2006; Roper, 2007). Other mechanisms include proton exchange, a stimulus-gated $Ca^{++}$
179 channel and the direct entry through an $H^+$ channel that has not been identified (Di Lorenzo *et al.*,
180 2009). A sour taste can also be induced by the passage of electric current through the tongue that
181 probably generates hydrogen ions from the hydrolysis of acid or water (Damodaran & Parkin, 2017;
182 Wong, 2018). In addition, undissociated acids play an important role in sour perception. For instance,
183 some weak organic acids that naturally occur in foods, such as citric, succinic, malic, or lactic acid,
184 are perceived to be more sour than hydrochloric acid at the same pH (Ley *et al.*, 2012). On the other
185 hand, other acid molecules (i.e., potassium acid oxalate or protocatechuic acid) exhibit both sour and
186 bitter tastes (Wong, 2018). The standard used to assess the sourness in food is citric acid (Liu *et al.*,
187 2020; Rojas *et al.*, 2022; Shiyan *et al.*, 2021; Yu *et al.*, 2018).

188 Saltiness is the sensation produced by some soluble salts, such as those with low molecular-weight,
189 mainly chlorides from sodium, potassium or calcium (Damodaran & Parkin, 2017; Wong, 2018).
190 NaCl is the only compound exhibiting an intense and clean (no after taste) salty taste and it is
191 consequently used as the saltiness standard (Liu *et al.*, 2020; Rojas *et al.*, 2022; Shiyan *et al.*, 2021;
192 Yu *et al.*, 2018). Potassium chloride can be considered a replacement for NaCl, however it can be
193 perceived as a sweet/bitter taste at low concentrations (Di Lorenzo *et al.*, 2009). In contrast, high
194 molecular-weight salts elicit bitter rather than salty taste, such as lithium chloride and ammonium
195 chloride. However, they are limited for human consumption due to safety and their offensive tastes,
196 respectively (Ley *et al.*, 2012; Wong, 2018).

The physiological function of the salty taste is to maintain the body's electrolyte balance. In taste buds, ion channels allow the passage of chemical species that trigger stimuli perceived as salty, and there is a relationship between the number of fungiform papillae and sensitivity to salty taste (Doty *et al.*, 2001). Apparently, the salty taste is related to the body's ability to detect sodium, thanks to the specific transduction mechanism of this cation, and its passage through the epithelial-sodium channel (ENaC) in the apical membrane of the receptor cells of taste. The epithelial-sodium channel is the mammalian $Na^+$ specific taste receptor. Most mammals have at least one type of salt taste receptor that is cation nonselective, apparently from the salty taste evoked by KCl and $NH_4Cl$ molecules. At the same time, high circulating aldosterone levels suggest aldosterone modulated epithelial cell membrane $Na^+$ transporters as candidate for salt taste receptors (DeSimone & Lyall, 2006). Moreover, one or more receptors, such as a variant of TRPV1 (TRPV1t), may be able to respond to various cations including $K^+$, $Ca^{2+}$, $NH_4^+$ and to $Na^+$ (DeSimone & Lyall, 2006; Rhyu *et al.*, 2021).

### 1.1.3. Non-basic tastes

Some compounds or combinations of compounds can produce tastes considered as non-basic or secondary tastes, such as astringent, chilling, cooling, heating and pungent (Damodaran & Parkin, 2017; Ley *et al.*, 2012; Wong, 2018). Moreover, other sensations have also been described as non-basic tastes, such as fattiness, or the definition of water as a tastant. Other characteristics of substances have led to the classifications of compounds as multitastes, taste modifiers or tasteless.

The definition of fattiness as a taste has been triggered by the transduction mechanisms that are sensitive to fatty acids in the TRCs membranes (Gilbertson *et al.*, 1997). The transduction mechanisms are associated with the inhibition of delayed rectifying $K^+$ channels and through the fatty acid CD36 (Di Lorenzo *et al.*, 2009). Evidence suggests that fatty acids (e.g. linoleic acid, oleic acid and stearic acid) could be considered as tastants and that their tastes are detectable without the need for other sensory cues such as texture, viscosity or smell (Di Lorenzo *et al.*, 2009). On the other hand, water has its own taste, even though it could be affected by temperature and easily affected by diluted compounds even at low concentration. Moreover, it could be considered as a tastant because of the role of water in eliciting compounds in TRCs and in taste nerves of some species (Di Lorenzo *et al.*, 2009). It has been suggested that an aquaporin, AQP5, a membrane channel, allows the water molecules to get into the cell by activating and regulating the volume of water through the anion channel (Di Lorenzo *et al.*, 2009). Moreover, when the mouth is rinsed after the application of a sweet taste blocker, water elicited a sweet aftertaste (Di Lorenzo *et al.*, 2009).

Multitaste is a complex sensation of tastes elicited by combining more than one basic taste at the same time (Rojas *et al.*, 2022). It is triggered by a variety of different compounds. Some examples of

multitaste compounds are the potassium acid oxalate and protocatechuic acid, which produce sour/bitter tastes (Wong, 2018), calcium phenolsulfonate (bitter/astringent tastes) and benzyl acetate (bitter/pungent tastes) (Dagan-Wiener *et al.*, 2019). Some compounds are able to alter and even block the taste of other compounds. $Na^+$ channel blockers reduce the saltiness of sodium chloride, thaumatin and adenosine monophosphate block bitterness, while lactisol proprionate blocks sweetness. On the other hand, some compounds increase the taste of others (taste enhancers); for example chlorogenic acid and cynarin enhance the sweetness (Di Lorenzo *et al.*, 2009). In contrast, some compounds have antagonist effects, that is, they tend to suppress the taste sensation of other compounds. This is what occurs with citric acid and sucrose tastants in lemonade. It is also possible to find synergistic effects, for instance the enhancement of umami taste by the addition of IMP or GMP to MSG (Di Lorenzo *et al.*, 2009).

The expression "tastelessness" is used to categorize molecules as lacking any particular taste. These are also classified as non-sweet, non-bitter, non-sour, non-salty or non-umami compounds (Rojas *et al.*, 2017; Rojas *et al.*, 2022). Some changes in the chemical structure of substances may modify their sweet taste to a bitter one or make them tasteless. For example, the saccharin sweetener becomes bitter by the introduction of a nitro group onto carbon five (5-nitrosaccharin), while the introduction of this group on the four-carbon position produces a sweet/bitter tastant (*p*-nitrosaccharin). On the other hand, the presence of the amino group produces a sweet/tasteless compound (6-aminosaccharin) or a tasteless molecule (5-aminosaccharin) (refer to Figure 1) (Rojas *et al.*, 2022). Interestingly, some tasteless compounds like miraculin and circulin act as taste modifiers, in particular, these compounds change the sense of sour in substances to sweet. In contrast, gymnemic acid, ziziphin and hodulcin block the sensation of sweetness (Di Lorenzo *et al.*, 2009).


**Figure 1 should be inserted around here**


## 1.2. Machine learning to uncover Structure-Taste Relationships

Studies of Quantitative Structure-Property Relationships (QSPRs) enhance the definition of mathematical relationships between molecular structures and specific properties of chemical compounds, such as taste. These approaches have played an important role in the evaluation and study how molecular features are related to the taste of chemical substances through the development of empirical data-driven models. QSPR models require molecular descriptors, which are numerical indices that encode the detailed chemical and structural information of molecules. They can be both experimental physicochemical properties of molecules and theoretical indices, which are calculated

through mathematical algorithms (Todeschini & Consonni, 2009). Molecular descriptors are used as independent variables in QSPR models. The relationships between descriptors and the property of interest (e.g., the taste of chemicals) are calculated by means of chemometrics and machine learning approaches.

The QSPR workflow starts with an appropriate description of the molecular structures and ends with the prediction of the behavior of the chemicals. This approach relies on the assumption that the molecular structure of a substance encodes the chemical features that are responsible for its physical, chemical, and biological behavior. If these features are correctly encoded into numerical descriptors, then QSPR strategy allows first to establish the empirical relationships between descriptors and the property of interest by means of statistical multivariate modeling, and subsequently infers the property of a new substance or untested chemical through the QSPR model.

There are several multivariate statistical methods to process molecular descriptors and achieve reliable estimates of chemical properties. Depending on the nature of the modeled property, classification and regression methods can be used to calculate models both for reproducing the known experimental data and predicting the unknown data for qualitative and quantitative responses, respectively. If chemicals belong to defined qualitative classes; for example, molecules labelled as positive or negative, then supervised classification models can be applied. Classification approaches define mathematical relationships between descriptors and classes and can thus be used to predict the class of new substances that are associated with unknown experimental class labels. If chemicals are associated with a quantitative response, regression methods are used to define the mathematical model that relates descriptors and the response to obtain quantitative predictions for new chemicals.

The two main operational steps in the development of QSPR models are the definition of their applicability domain and the implementation of proper validation protocols, as proposed by the OECD (Organization for Economic Co-operation and Development) in the framework of the five general principles for QSARs (Gramatica, 2007). These principles are used as the criteria to evaluate and accept QSPRs, especially for regulatory purposes, and state that each model should have: 1) a defined endpoint; 2) an unambiguous algorithm; 3) a defined domain of applicability; 4) appropriate measures of goodness-of-fit, robustness and predictivity; 5) a mechanistic interpretation, if possible. The Applicability Domain (AD) of a QSPR model is the chemical space where predictions can be considered as reliable (Mathea *et al.*, 2016; Sahigara *et al.*, 2012). If the properties of a new untested molecule are predicted through QSPRs, then it is considered to share the same mechanisms and/or modes of action as the molecules used to build the model provided that it is structurally similar to the training molecules and falls inside the AD. In this case, the properties of predicted chemical are considered as interpolated by the model and its predicted properties can be assumed to be reliable. In

297  contrast, the predictions for molecules falling outside the AD can be considered as model
298  extrapolations, and consequently they are considered to be unreliable.

299  Moreover, the attention to effective and reliable estimates through predictive models has a crucial
300  role in the QSPR workflow. When supervised qualitative (classification) or quantitative (regression)
301  approaches are used to establish structure-property relationships, the primary goal of the process is to
302  achieve reliable models that are able to correctly predict the properties of new untested molecules.
303  QSPR modeling could also have explanatory purposes, that is, allowing the interpretation of the
304  relationship between descriptors and the modeled property to deepen the knowledge about the specific
305  problem in analysis. In both cases, validation protocols for the assessment of the predictive ability of
306  models and the reliability of the established relationships should be always applied (Oliveri, 2017;
307  Wold & Eriksson, 1995). This step is necessary also to avoid overfitted models, that is, models in
308  which mathematical relationships accurately predict properties for the training compounds, but not
309  for new untested substances.

310  The predictive abilities of the models are usually evaluated by splitting the available compounds into
311  training and test sets. Training compounds are used to establish the mathematical model, which is
312  then used to predict the responses of the chemicals included in the test set. Finally, the agreement
313  between of experimental and predicted responses for the test substances is evaluated to assess the
314  model's predictive ability. Several validation protocols exist and the usage of a particular one usually
315  depends on how many chemicals are available for model development. A general requirement is that
316  molecules in the test chemical space should be reasonably similar to that of the training space.
317  However, large degrees of similarities could produce an excessively optimistic evaluation of a
318  model's predictive ability. For this reason, when dealing with classification models, it is preferable
319  to keep the class balance equal in the training and test sets; that is, the same distribution of chemicals
320  in the modeled classes should be preserved in both sets.

321  **1.2.1. Classification approaches**
322  In the framework of machine learning applied to QSPR, classification methods are fundamental
323  techniques aimed at finding mathematical relationships that recognize the class membership of
324  molecules on the basis of a set of molecular descriptors. Once a classification model has been trained,
325  the membership of unknown chemicals to one of the defined classes can be predicted. Thus, for
326  discrete molecular properties, like qualitative properties distinguishing between different tastes, a
327  general representation of classification models is the following:

328

329 
$$C = f\left(x_1, x_2, \ldots, x_p\right)$$
(1)

10

330

where C is the class, $x_1$, ..., $x_p$ are $p$ (number) of molecular descriptors, and $f$ is a function representing the relationship between the class and the descriptors.

Several classification methods have been proposed in the last decades, with different characteristics, advantages and limitations (Lavine & Rayens, 2009). A preliminary distinction among classification methods can be defined on the basis of the mathematical form of the decision boundary: linear methods calculate the best linear boundary for class discrimination, while non-linear methods discriminate classes by non-linear boundaries.

Another important difference can be made between discriminant (pure classification) and class-modeling methods. Discriminant methods divide the whole chemical space defined by the molecular descriptors in as many regions as the number of the modeled classes. Thus, each compound is assigned the class corresponding to the region of the chemical space where it falls. On the other hand, class-modeling methods (also known as one-class classifiers) define the boundary to separate a specific class from the rest of the chemical space. Thus, a target class is modeled independently of the others; compounds fitting the class model are considered members of the class, while chemicals that are outside the class space are classified as non-members of the target class.

Among classification methods, Discriminant Analysis (DA) is the most widely used (Hand, 1997; McLachlan, 1992). DA finds the directions in the multivariate space that maximizes the ratio of the between-class to within-class variances; these are called discriminant functions and from a mathematical point of view, these directions are linear combinations of the original variables. Depending on the choice of the class-covariance representation, two different discriminant methods can be distinguished: Quadratic Discriminant Analysis (QDA) and Linear Discriminant Analysis (LDA), which define quadratic and linear boundaries between classes, respectively. One major drawback for DA is that it cannot be applied to datasets with the number of samples lower than the dimension of the measurement space. However, to overcome this limitation, DA can be combined with methods for dimensionality reduction, such as variable selection approaches or principal component analysis (PCA).

Another option to deal with highly dimensional spaces is the application of Partial Least Squares Discriminant Analysis (PLSDA) (Barker & Rayens, 2003; Brereton & Lloyd, 2014). PLSDA benefits from the properties of PLS (Partial Least Squares) regression, since it searches for the latent variables, that is, the directions of maximum covariance with the response to be modeled. The difference from PLS is that the response encodes class membership with binary codes and class thresholds have to be defined to predict samples in one of the modeled classes.

363   Unlike DA and PLSDA, which are discriminant classifiers, the Soft Independent Modeling of Class

364   Analogy (SIMCA) method is one of the most useful and popular class-modeling approaches. It is

365   based on PCA carried out on the samples of the target class. To predict the class of test samples, the

366   sample distances from the class PCA model are calculated on the basis of normalized Q residuals and

367   Hotelling's $T^2$ values, which measure how well each sample conforms to the model. Only samples

368   with distances lower than a defined threshold are classified into the class space.

369   Another class modeling approach consists of the calculation of Potential Functions (PFs), where the

370   assignment of a new sample to the target class is based on the cumulative potential of the class, which

371   is calculated as the sum of the individual potentials of the target class samples in the point of the

372   chemical space where the new sample is projected. The shape of the potential depends on the choice

373   of the type of potential function (kernel) and smoothing parameter (Brereton, 2011).

374   Tree-based algorithms exploit different classification approaches. They recursively divide data into

375   smaller subgroups, which contain samples belonging to as few classes as possible. In each split, the

376   partition is achieved by maximizing the purity of the new subsets. The final classification model

377   consists of a collection of nodes that define the classification rule. One of the most common tree-

378   based approaches is the Classification and Regression Tree (CART), which selects the variables that

379   provide the purest subsets of samples in each node (Breiman *et al.*, 1984). The Random Forest (RF)

380   method represents a subsequent development of tree-based approaches (Breiman, 2001). It is a meta-

381   classifier based on an ensemble of classification trees, each trained on various subsamples of the

382   training set, which are built by bootstrapping. The prediction is then obtained by majority vote among

383   the classifications provided by the trees of the forest.

384   Another approach, based on the ensemble of models, is AdaBoost (Adaptive Boosting), where

385   predictions provided by many "weak" classifiers are pooled to produce a better classification.

386   Predictions are combined through an adaptive iterative algorithm that exploits the weighted majority

387   voting (Freund & Schapire, 1997). Besides the original boosting method, other approaches have been

388   proposed and applied for the prediction of molecular taste, especially when dealing with big datasets,

389   such as XGBoost (eXtreme Gradient Boosting) (Chen & Guestrin, 2016). This is again a classification

390   algorithm that uses sequential iterations, where decision trees are combined to increase classification

391   accuracy.

392   Often QSPRs exploit similarity-based classification, since compounds with similar molecular

393   structures are expected to have similar properties. These methods calculate distance measures to

394   provide a classification in terms of similarity among samples. The most known approach in this

395   framework is the *k*-Nearest Neighbors (*k*NN) classifier: it classifies a sample according to the most

396   frequent class of its *k* most similar training samples (Kowalski & Bender, 1972). The N3 (*N*-Nearest

Neighbors) approach is an evolution of $k$NN, which uses locally-weighted information to classify new samples. The Binned Nearest Neighbors (BNN) method is similar to $k$NN, but the prediction is based on a flexible number of neighbors (Todeschini *et al.*, 2015a).

Another classification approach, which is relatively frequent in QSPR applications, is the Support Vector Machine (SVM) method (Vapnik, 1998). It defines the boundary between two classes by maximizing the distance between the support vectors and the decision boundary, where support vectors are those training samples located in the proximity of the class border. Moreover, SVM can use non-linear kernel functions for defining non-linear decision boundaries.

To visually exemplify the different ways classification methods can define boundaries between classes, a dataset of 324 chemicals was generated from the ChemTastesDB database (Rojas *et al.*, 2022), including 61 chemicals labelled as sweeteners and 263 as bitterants. Their molecular structures were encoded through the binary molecular access system (MACCS) keys (Durant *et al.*, 2002). The chemical space was represented by the first two t-Distributed Stochastic Neighbor Embedding (t-SNE) dimensions (van der Maaten & Hinton, 2008), calculated by using the Jaccard-Tanimoto metric as the distance measure (Todeschini *et al.*, 2015b). Finally, different classification approaches were calculated to show how the class boundaries can vary in a 2D space according to the adopted method (Figure 2). SIMCA and PFs, which are class-modeling approaches, define a boundary around the target class (sweet class in this example), while the discriminant methods, for instance, LDA, QDA and PLSDA, divide the entire chemical space into two sub-spaces, each associated with one of the two modeled classes, with linear or non-linear boundaries, depending on the adopted classification algorithm.


**Figure 2 should be inserted around here**


### 1.2.2. Classification measures

QSPR models must be assessed through measures of goodness-of-fit and goodness-of-prediction. In this framework, several indices can be used to evaluate the quality of models, which are based on the number of misclassifications (molecules assigned to the wrong class) (Ballabio *et al.*, 2018). Classification metrics are derived from the confusion matrix, which is a square matrix with dimensions $G \times G$, where $G$ is the number of modeled classes. Each entry $c_{gk}$ of this matrix represents the number of samples belonging to class $g$ and assigned to class $k$. Consequently, the diagonal elements $c_{gg}$ denote the counts of the correctly classified samples while the off-diagonal elements represent those erroneously classified. In the simplest binary case where two classes (positive and

13

430  negative) are modeled, the confusion matrix is a $2 \times 2$ numerical table with four entries labelled as

431  follows: true positive and true negative (TP and TN, the number of positive and negative samples

432  correctly classified, respectively), false positive (FP, the number of negative samples classified as

433  positive) and false negative (FN, the number of positive samples classified as negative).

434  The most common classification measures derived from the confusion matrix are *sensitivity* ($Sn_g$),

435  *precision* ($Pr_g$), *specificity* ($Sp_g$), as well as their combination, such as the *F-score* ($F_g$) (also known

436  as the *$F_1$-score* or *F-measure*). These indices are associated to each modeled *g*-th class and defined

437  as:

$$Sn_g = \frac{c_{gg}}{n_g} \qquad\qquad Pr_g = \frac{c_{gg}}{n'_g}$$

438

$$Sp_g = \frac{\sum_{\substack{k=1 \\ k \neq g}}^{G} \left( n_k - c_{kg} \right)}{n - n_g} \qquad\qquad F_g = 2 \cdot \frac{Sn_g \cdot Pr_g}{Sn_g + Pr_g}$$
(2)

439

440  where $n_g$ is the number of samples of the *g*-th class, $n'_g$ is the number of samples that are classified

441  in the *g*-th class and *n* is the total number of samples. Higher values of sensitivity, specificity and

442  precision are associated with better class discrimination.

443  Beside measures assigned to each class, global classification indices have been proposed to provide

444  an overall assessment of the discrimination ability of classifiers. The *Non-Error Rate* (*NER*, also

445  called *balanced accuracy* or *recall*) corresponds to the arithmetic mean of class sensitivities:

446
$$NER = \frac{\sum_{g=1}^{G} Sn_g}{G}$$
(3)

447

448  while *accuracy* corresponds to the fraction of correctly classified samples:

449
$$ACC = \frac{\sum_{g=1}^{G} c_{gg}}{n}$$
(4)

450

451  Note that accuracy is considered a biased estimate when classes are unbalanced, that is, samples are

452  distributed in classes with significantly different frequencies.

453  Alternatively, classification performance can also be evaluated through the Matthew Correlation

454  Coefficient (*MCC*), which ranges between -1 and 1 and has originally been defined to assess binary

455  classification tasks:

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FN) \cdot (TP + FP) \cdot (TN + FP) \cdot (TN + FN)}} \qquad (5)$$

Another way to assess discrimination capability of classification models is through ROC (Receiver Operating Characteristics) curves. These are graphic plots of *sensitivity* and 1 - *specificity* (also known as False Positive Rate, FPR) for a classification system when its discrimination threshold is changed. For each threshold value, the corresponding TPR and FPR values are calculated. The optimal classifier will provide a full ROC curve, while a random classification rule would give a line along the diagonal of the ROC space. To quantitatively compare classification models trough ROC curves, a common approach is to calculate the area under the curve (*AUC*), also known as *AUROC* or *ROC-AUC*.

## 2. Classification models for taste prediction

In this section, ligand-based (LB) classifiers for taste prediction are described. The classifiers were retrieved from 52 published studies, which were found through critical screening of the Web of Science citation indexing service. To the best of our knowledge, the first work published on this topic was by Lemont B. Kier in 1980. In addition to the models described below, there are several multicriteria reviews that are focused on QSAR-based prediction of tastes by means of diverse classification-based machine learning approaches (De León *et al.*, 2021; Malavolta *et al.*, 2022; Rojas *et al.*, 2016a; Spillane *et al.*, 1996; Walters, 2006). Ligand-based models are presented on the basis of the basic tastes to be predicted.

### 2.1. Sweet and bitter tastants

The discrimination between sweet and bitter compounds has probably been the most important task in quantitative structure-taste relationship studies. Twelve ligand-based (LB) models for the discrimination of these two tastes are summarized in Table 1.

**Table 1 should be inserted around here**

Earlier studies had been focused on the use of simple modeling approaches, such as discriminant analysis. In 1980, Kier (Kier, 1980) performed a two-variable linear discriminant analysis (LDA) to discriminate sweet and bitter aldoximes taken from the data published by Acton and Stone (Acton & Stone, 1976). For each class, 10 tastants were selected on the basis of the largest percentage of the taste and the most potent taste response. Each molecule was represented by two connectivity indices

15

487 named $^1\chi$ and $^4\chi_p$. The classifier was used to predict the taste of nine external molecules, achieving

488 seven correct predictions, one incorrect prediction and one tastant labelled as ambiguous. After this

489 pioneering work, Takahashi and Miyashita's group (Miyashita *et al.*, 1986a; Takahashi *et al.*, 1984;

490 Takahashi *et al.*, 1982) developed new models, based on LDA and SIMCA. In the first study

491 (Takahashi *et al.*, 1982), three molecular descriptors were used to correctly classify 22 perillartines

492 (11 in each class) through a LDA classifier. In a subsequent study (Takahashi *et al.*, 1984), a test set

493 of nine compounds (five sweet and four bitter) retrieved from Acton's dataset (Acton & Stone, 1976)

494 was included. Two LDA classifiers, one based on three descriptors and one with just two descriptors,

495 were developed, achieving similar performances on both training and test sets. In the last study,

496 Miyashita (Miyashita *et al.*, 1986a) used 70 sweet and 21 bitter aspartyl dipeptides (L−Asp−NH−R)

497 to calibrate a five-variable SIMCA model.

498 Drew (Drew *et al.*, 1998) used a dataset of 21 sweeteners, 20 sweet/bitter and 9 bitter mono- and di-

499 substituted sodium sulfamates, which were properly optimized by the semiempirical PM3 method to

500 calculate 11 molecular descriptors. Then, they performed a discriminant analysis (DA), which was

501 able to perfectly discriminate all the compounds. In addition, a cluster analysis was carried out in the

502 space of the first two principal components, where a linear separation could be found only between

503 sweet and sweet/bitter tastants. Two years later, Spillane (Spillane *et al.*, 2002) synthetized and

504 characterized 23 meta-phenylsulfamate derivatives. The equilibrium geometry of the tastants were

505 obtained by means of the AM1 semi-empirical method, in such a way as to calculate diverse

506 descriptors, which were used to calibrate a discriminant plot, an LDA and a quadratic discriminant

507 analysis (QDA). The first model was obtained by plotting the values of length ($x$, Å) against the

508 volume $V_{CPK}$ ($xyz$, Å$^3$). The best LDA classifier was obtained with the $x$, width ($z$, Å), aqueous

509 solvation energy ($E_{solv}$) and HOMO descriptors; while the best QDA model used the $x$, $z$, $E_{solv}$ and

510 LUMO descriptors. Among these models, the QDA exhibited the best performance in terms of the

511 *NER*. In a further analysis, these three models were used for predicting the taste of 9 unsynthesized

512 meta-compounds.

513 Other models to discriminate sweetness and bitterness were based on the *k*-Nearest Neighbors (*k*NN)

514 approach. The first model was proposed by Takahashi (Takahashi *et al.*, 1982) to classify 22

515 perillartines. The *k*NN method misclassified only two bitter molecules in the entire dataset. Several

516 years later, *k*NN was used by Rojas (Rojas *et al.*, 2016c) with 508 curated and filtered tastants (427

517 sweet and 81 bitter), which were split into training (356 tastants) and test sets (152 molecules).

518 Molecules were represented by means of 3,763 conformation-independent Dragon molecular

519 descriptors (Kode srl., 2018), which were initially analyzed by means of the V-WSP unsupervised

variable reduction approach (Ballabio *et al.*, 2014). Then, the training set was used for model development using the 5-fold cross-validation approach to determine the optimal $k$ value during the genetic algorithms-variable subset selection (GAs-VSS). A four-descriptor model was selected as optimal, with a balanced performance in prediction ($NER = 0.789$, $Sn_{sweet} = 0.953$ and $Sn_{bitter} = 0.625$). In addition, the applicability domain (AD) of the model was calculated. In a further analysis, the sweeteners database was used to perform a quantitative structure-property relationship (QSPR) for predicting the relative sweetness (RS) or sweetness potency (Sw) of the sweeteners (Rojas *et al.*, 2016b).

Starting from 2017, the random forest (RF) classifier started to be applied to discriminate sweet and bitter tastants. Chéron (Chéron *et al.*, 2017) merged 316 sweeteners from SweetenersDB and 680 bitterants from the BitterDB (Wiener *et al.*, 2012), which were represented by 244 conformation-independent Dragon descriptors (Kode srl., 2018). One hundred trees with a tree depth of five molecular descriptors and the Gini splitting criterion were set up during the calibration of the RF classifier, which exhibited good performance on the test set, constituted by the 20% of tastants ($NER = 0.914$ and $MCC = 0.848$). In a further analysis, the model was used to identify 4,585 natural molecules of the SuperNatural II database (Banerjee *et al.*, 2015) as potential sweet agents and their relative sweetness was predicted by means of a support vector regression (SVR). One year later, Banerjee and Preissner (Banerjee & Preissner, 2018) calibrated a RF model, named BitterSweetForest, for 517 sweeteners from SuperSweet (Ahmed *et al.*, 2011) and 685 bitterants from the BitterDB (Wiener *et al.*, 2012). Compounds were represented by means of the extended connectivity fingerprints ECFP4 (Morgan, 1965; Rogers & Hahn, 2010) calculated using RDKit. The best model achieved good predictive ability of the 241 test set molecules ($AUC = 0.98$, *F-score* = 0.92, $ACC = 0.967$, and *Cohen's Kappa* = 0.92). In addition, the BitterSweetForest model was used to virtually screen the SuperNatural II (Banerjee *et al.*, 2015) and DrugBank databases. Goel (Goel *et al.*, 2021) developed a dataset of 1,179 sweeteners and 743 bitterants (retrieved from the BitterSweet database (Tuwani *et al.*, 2019)) and used the recursive feature elimination approach to identify eight descriptors from 1,613 conformation-independent Mordred descriptors (Moriwaki *et al.*, 2018). The best RF classifier exhibited good prediction on the test set (20% molecules), with $NER = 0.855$, $ACC = 0.865$ and $MCC = 0.785$. In addition, 478 structurally diverse sweeteners (334 in the training set and 144 in the test set) were used to predict the relative sweetness (log RS) by means of a 3D regression-based RF model, which was then submitted to a molecular docking simulation to calculate the binding conformation and associated free binding energy with the T1R2/T1R3 receptor. In a subsequent step, compounds from the Universal Natural Products Database (UNPD) (Gu *et al.*,

2013) were virtually screened following the above mentioned workflow, which was coupled with toxicity scaffold analysis as well.

Recently, other advanced classifiers were used for sweetness and bitterness discrimination. In 2022, Bo (Bo *et al.*, 2022) curated a dataset of 797 bitterants and 1,249 sweeteners retrieved from the BitterDB (Dagan-Wiener *et al.*, 2019), SuperSweet (Ahmed *et al.*, 2011) and FlavorDB. The 2D RDKit molecular descriptors and fingerprints were used to calibrate multilayer perceptron (MLP) models, while 2D-RGB color images ($32 \times 32$ pixels) were used to train convolutional neural networks (CNN). Among the three models, the best one (BitterSweetMLP-Fingerprint) was obtained with 17 fingerprints (selected by means of the PCA using oblique rotation), with good performance for predicting the 409 test set tastants ($NER = 0.880$, $AUC = 0.950$, $ACC = 0.880$, and $MCC = 0.750$). Molecular charges and their surface interaction descriptors were retained since they were relevant for classifying sweeteners/bitterants. Maroni (Maroni *et al.*, 2022) calibrated a gradient boosting machine model (LightGBM implementation), along with other well-known classifiers: $k$NN, RF, logistic regression (LR) and multilayer perceptron (MLP). These authors filtered and curated a database of 2,195 tastants, which were represented by 1,402 conformation-independent features calculated in the RDkit, Pybel (O'Boyle *et al.*, 2008) and Mordred (Moriwaki *et al.*, 2018). A sequential descriptor selection combined with hierarchical clustering in the descriptor's Spearman rank-order index was used. The GBM classifier was optimal with good results in calibration ($NER = 0.893$, $AUC = 0.950$ and $F\text{-}score = 0.883$). Furthermore, the SHapley Additive exPlanations (SHAP) allowed the identification of the most suitable descriptors.

**2.2. Sweet and non-sweet tastants**

Several classification-based machine learning models have been built to discriminate between sweet and the non-sweet molecules, as well as to use them in order to predict and synthesize novel sweeteners. Nineteen ligand-based classifiers for sweet taste predictions are summarized in Table 2.

**Table 2 should be inserted around here**

As presented in the sweet/bitter section, the sulfamate sweetness prediction was based on models using biplot discriminant analysis (DA) published mainly by Spillane's research group. In the first application, Spillane and McGlinchey (Spillane & McGlinchey, 1981) used the length ($x$, Å) and volume $V_{CPK}$ (xyz, Å$^3$) descriptors to construct a DA-biplot for the discrimination of 47 sweet and non-sweet carbosulfamate ($RNHSO_3^-$) derivatives. In a second study, Spillane and Sheahan (Spillane & Sheahan, 1989) again used the $x$ and $V_{CPK}$ descriptors to classify 17 carbosulfamates. In a

586 subsequent DA-plot application (Spillane *et al.*, 1993), the $V_{CPK}$ and $\Sigma \sigma$ descriptors were used for
587 40 synthesized ring disubstituted phenylsulfamates as sodium salts (no classification performances
588 were reported for this discriminant plot). Between 1983 and 2009, the same research group developed
589 six models based on linear discriminant analysis (LDA) and four models based on quadratic
590 discriminant analysis (QDA). In the first LDA application, 33 sweet and non-sweet heterosulfamates
591 were used (Spillane *et al.*, 1983). Twenty molecules were retrieved from the Acton's database (Acton
592 & Stone, 1976), while another 13 tastants were synthesized and evaluated by the authors for taste
593 sensation. The best model was composed of the length ($x$, Å), width ($z$, Å) and the first-order valence
594 connectivity index ($^{1}\chi^{v}$) descriptors. In a subsequent study, Spillane and Sheahan (Spillane &
595 Sheahan, 1989) used the same pool of descriptors to calibrate a LDA model for other 23
596 heterosulfamates and an extended dataset of 56 heterosulfamates.
597 Starting from 2000, the classification and regression tree (CART) method was also used for sweetness
598 prediction. Spillane (Spillane *et al.*, 2000) augmented previous datasets in order to include 101
599 heterosulfamate sodium salts (32 were synthesized for this study). The datasets contained 20
600 sweeteners and 81 non-sweet derivatives. LDA and QDA models were calibrated with four molecular
601 descriptors ($x$, $y$, $z$ and and $^{1}\chi^{v}$), while with CART, three features were used ($x$, $y$ and $^{1}\chi^{v}$). Among
602 these models, the QDA classifier showed the best performance. Three years later, the dataset was
603 further augmented by including newly synthesized compounds (15 sweet and 16 non-sweet) (Spillane
604 *et al.*, 2003). In this case, CART provided better results than LDA and QDA using the $y$, $z$, $V_{CPK}$ and
605 LUMO descriptors. In 2005, Kelly (Kelly *et al.*, 2005) merged 63 sweeteners available in the
606 literature with 19 cyclamate derivatives that were synthesized and tasted in this work. The sweetness
607 value was used to define three classes of the predominant tastes: non-sweet (0 to 39), sweet/non-
608 sweet (40 to 60) and sweet (61 to 100). The dataset was randomly split (maintaining the class
609 proportion) into a training set and a test set of 75 molecules and 8 molecules, respectively. In this
610 work, an external validation was used for the first time for the sweet/non-sweet discrimination. The
611 LDA and QDA models exhibited poor predictive ability, while CART based on six descriptors ($x$,
612 HOMO, LUMO, $E_{solv}$, $V_{Spartan}$ and σ) exhibited acceptable prediction performance.
613 One year later, Spillane (Spillane *et al.*, 2006) developed three CART models to study a dataset of 82
614 tastants (42 newly synthesized disubstituted phenylsulfamates). The best classifier used 70 molecules
615 in the training set and 12 test set compounds randomly selected. Molecules in the test set were only
616 the newly synthesized non-sweet (11 compounds) and sweet/non-sweet (1 compound), while the four
617 newly synthesized sweeteners were placed in the training set. This model used seven descriptors and
618 provided good prediction ability. Finally, 28 five-membered aromatic ring thiazolyl-, benzothiazolyl-,

and thiadiazolylsulfamates were synthesized and merged together with 30 well-known heterocyclic sulfamates to create a database (Spillane *et al.*, 2009). Compounds were grouped into three classes according to the predominant taste: sweet, non-sweet and sweet/non-sweet. LDA and QDA were initially used considering all the molecules as training chemicals. Then, the authors calibrated two CART models by randomly splitting the dataset into a training set (48 tastants) and a test set (10 molecules). Between these two models, the best CART classifier used six descriptors and exhibited a moderate performance when applied to the test set.

In another two studies, Miyashita's group (Miyashita *et al.*, 1986b; Okuyama *et al.*, 1988) also used sulfamate derivatives to calibrate structure-taste relationships based on the SIMCA classifier. In the first work, 14 sweet and 36 non-sweet carbosulfamates described by molar refractivity (MR), five geometrical STERIMOL features and the Taft's $\sigma^*$ descriptor were used (Miyashita *et al.*, 1986b). The SIMCA model correctly predicted 13 sweet and 24 non-sweet molecules. In addition, a set of alkyl groups were proposed as potential substituents, from which six alkylsulfamates were predicted as potential sweeteners. Among these compounds, one was synthesized and exhibited a relative sweetness of three times greater with respect to sucrose. Two years later, the same authors used 25 acyclic and 20 cyclic carbosulfamates represented by different graph theoretical invariants (Okuyama *et al.*, 1988). In addition, the acyclic sulfamates were also represented by the weighted path numbers for the rooted atom (path length from 1 to 8) and counts of self-returning walks for the rooted atom (number of steps from 2 to 13), while the atomic path numbers for the rooted atom (path length from 1 to 8) and the counts of self-returning walks for the rooted atom were also computed for cyclic sulfamates. In both cases, the SIMCA model achieved similar performance for the acyclic carbosulfamates and the cyclic derivatives.

The first *k*NN model for the discrimination between sweet and non-sweet tastants was published in 2016 (Rojas *et al.*, 2016c). A nine-descriptor *k*NN model provided the best discrimination between 433 sweet and 133 tasteless curated molecules, with similar performances for training (*NER* = 0.838) and test sets (30% of compounds), *NER* = 0.752. One year later, the same research group (Rojas *et al.*, 2017) developed an expert system that integrated unsupervised and supervised machine learning approaches. To this end, a database of 435 sweet and 214 non-sweet (bitter and tasteless) molecules were represented by means of 875 conformation-independent descriptors (Todeschini & Consonni, 2009) and extended connectivity fingerprints (ECFPs) (Rogers & Hahn, 2010), calculated by the Dragon software (Kode srl., 2018). Similarity analysis was based on the ECFPs and multidimensional scaling (MDS), while the supervised classification was carried out with the consensus predictions provided by *N*-Nearest Neighbors (N3) and partial least squares discriminant analysis (PLSDA), with good predictive accuracy on the test chemicals (*NER* = 0.848, non-assigned = 19.3%). A new

consensus model was published in 2019 by Zheng (Zheng *et al.*, 2019) for a curated database of 530 sweet and 850 non-sweet molecules, which were represented by four types of ECFPs (Rogers & Hahn, 2010): 1024bit-ECFP4, 2048bit-ECFP4, 1024bit-ECFP6 and 2048bit-ECFP6. They used the *k*NN classifier, along with support vector machine (SVM), random forest (RF), gradient boosting machine (GBM) and deep neuron network (DNN) approaches to developed 1,312 individual models, as well as 96 averaged classification models. As a result, four consensus models were constructed (CM01 - CM04), and the best one (using 19 best individual models) was selected to construct the e-Sweet model. This model achieved good results in predicting the 221 test set compounds (*NER* = 0.900, *F-score* = 0.878 and *MCC* = 0.807). In a further step, a consensus regression was developed to predict the relative sweetness of 352 sweeteners.

In addition, Tuwani (Tuwani *et al.*, 2019) calibrated diverse models based on RF, ridge logistic regression and AdaBoost for the classification of sweet/non-sweet and bitter/non-bitter molecules (refer to bitter and non-bitter section). These models were named BitterSweet. For sweetness prediction, a dataset of 1,205 sweeteners and 1,171 non-sweeteners were represented by means of diverse molecular descriptors calculated in Dragon (Kode srl., 2018), Canvas (Schrödinger LLC, 2017) and ChemoPy (Cao *et al.*, 2013). The best model in terms of classification accuracy for the test set (7% of molecules) used the 2D/3D Dragon descriptors reduced by means of the Boruta algorithm and subsequently coupled with AB machine learning: *NER* = 0.834, *AUC* = 0.883 and *F-score* = 0.856.

Two years later, Fritz (Fritz *et al.*, 2021) developed the VirtualTaste prediction platform for predicting the sweet taste of molecules based on the RF classifier (VirtualSweet model). The database included 2,011 sweet and non-sweet (bitter and tasteless) molecules that were curated and standardized from the SuperSweet database (Ahmed *et al.*, 2011) and from their previous BitterSweetForest database (Banerjee & Preissner, 2018). Molecules were represented by MACCS keys (Durant *et al.*, 2002) and Morgan molecular fingerprints (Morgan, 1965; Rogers & Hahn, 2010). The RF model achieved good external prediction on the 403 test set tastants (*NER* = 0.893, *AUC* = 0.951, *F-score* = 0.888 and *ACC* = 0.893). Furthermore, the VirtualSweet model was used to virtually screen molecules from the DrugBank database and from the SuperNatural II database (Banerjee *et al.*, 2015). One year later, Yang (Yang *et al.*, 2022) used the RF and the XGBoost classifiers, along with other approaches, to calibrate diverse models for a database named Taste DB. However, this name was previously proposed by Ruddigkeit and Reymond (Ruddigkeit & Reymond, 2014). This dataset contained six families of compounds: natural (973 sweeteners and 687 non-sweeteners), artificial (402 positive and 798 negative), carbohydrate (220 sweet and 238 non-sweet), non-carbohydrate (1,155 positive and 1,476 negative), nutritive (226 sweet and 268 non-sweet) and non-nutritive (1,149 positive and 1,464

negative). For validation purposes, the datasets were divided into training and test set in a proportion of 8:2. The best artificial sweeteners model (in terms of accuracy for the test set prediction) used the RF classifier and MACCS structural keys ($NER = 0.920$ and $AUC = 0.971$), while for the carbohydrate family of compounds, the XGBoost approach and Atom pairs descriptors ($NER = 0.926$ and $AUC = 0.974$) were used. The remaining four models were developed by means of the XGBoost classifier and MOE2d descriptors with the following performances for the test set: 1) natural molecules ($NER = 0.841$ and $AUC = 0.920$); 2) non-carbohydrate compounds ($NER = 0.867$ and $AUC = 0.947$); 3) nutritive sweeteners ($NER = 0.876$ and $AUC = 0.956$); and 4) non-nutritive molecules ($NER = 0.889$ and $AUC = 0.961$). In further analysis, these authors developed regression models to predict sweetness potency (log Sw).

More recently, Bo (Bo *et al.*, 2022) developed diverse quantitative structure-taste relationships based on MLP and CNN deep learning classifiers, following the same workflow as previously described in the sweet/bitter section. In this case, the dataset contained 1,119 sweeteners and 1,101 non-sweeteners (tasteless and bitter). The best two models, in terms of their predictive ability, are the SweetMLP-Fingerprint ($NER = 0.900$, $AUC = 0.940$, $ACC = 0.880$ and $MCC = 0.800$) and SweetCNN ($NER = 0.850$, $AUC = 0.900$, $ACC = 0.840$ and $MCC = 0.660$). These models were used to predict the taste of 902 tastants of the bitter data set. Lee (Lee *et al.*, 2022) used a fully connected network (FCN), along the RF, XGB and LGBM classifiers, to propose the soft-vote ensemble approach. The curated dataset included 1,237 sweeteners and 1,054 non-sweeteners retrieved from the BitterSweet database (Tuwani *et al.*, 2019), which were represented by means of eight 2D fingerprints and diverse molecular descriptors. Among the 44 different models, the best models were LGBM applied to layered fingerprints and alvaDesc descriptors (Mauri & Bertola, 2022). These two models were used to assemble the BoostSweet model for sweetness prediction by means of the soft-vote method that averages the prediction of each model. The BoostSweet classifiers achieved good performance for the test set (211 sweeteners and 248 non-sweeteners): $NER = 0.899$, $AUC = 0.961$ and $F\text{-}score = 0.907$.

**2.3. Bitter and non-bitter tastants**

The prediction of sweetness has been the predominant goal for research in the computational taste framework, probably because bitterness was usually linked to toxic compounds (as described for the alkaloids). However, in the last few years, models that predict bitterness have received considerably more attention due to the use of bitterants in several applications, particularly in food and pharmaceutical industries. In contrast to the sweet/bitter models, where the main purpose was sweetness prediction, comprehensive classification models for bitterness prediction are focused on

720     discriminating bitter from non-bitter tastants. The 14 ligand-based models found to date in the
721     literature are summarized in Table 3.

722

723                   **Table 3 should be inserted around here**

724

725     In 2006 Rodgers (Rodgers *et al.*, 2006) used the Naïve Bayes (NB) classifier for the bitterness
726     prediction of small molecules. The curated dataset was composed of 649 bitterant taken from
727     scientific literature and patents, and 13,530 hypothetical non-bitter molecules randomly selected from
728     the MDL Drug Data Repository (MDDR). All the compounds were represented by MOLPRINT 2D
729     circular fingerprints (aka Atom Environments) (Bender *et al.*, 2004), which were subjected to a
730     variable subset selection. Ten years later, Huang released the first online tool, namely BitterX (Huang
731     *et al.*, 2016), for bitterness prediction based on support vector machine (SVM) classifiers. Data of
732     bitterants and bitterant-TAS2R interactions were retrieved from the PubMed (Sayers *et al.*, 2021) and
733     BitterDB (Wiener *et al.*, 2012) databases. In this work, a ligand-based model and a receptor-based
734     model were developed. In both cases, datasets were randomly split into training (80%) and test sets
735     (20%) three times to avoid bias in the data splitting, while genetic algorithms (GAs) were used for
736     the supervised descriptor selection. The ligand-based model was developed from a database of 539
737     bitterant and 539 non-bitter molecules and 46 physicochemical descriptors. The mean accuracy and
738     the area under the curve used in the prediction of the three models were $ACC = 0.915$ and $AUC =$
739     $0.950$. On the other hand, the TAS2R receptor recognition model used 260 bitterants and 260 non-
740     bitter molecules (negative), and 20 physicochemical and 15 receptor descriptors with slightly lower
741     prediction quality ($ACC = 0.798$ and $AUC = 0.823$).
742     RF classifiers were also used for bitterness prediction (Fritz *et al.*, 2021; Tuwani *et al.*, 2019). Tuwani
743     published the BitterSweet model (Tuwani *et al.*, 2019) for the classification of bitterants, in which
744     they followed the same workflow as presented in the sweet/non-sweet section. The RF classifier,
745     coupled with PCA reduction of ChemoPy descriptors, achieved a higher non-error rate in prediction
746     for the test set (154 molecules): $NER = 0.819$, $AUC = 0.880$ and $F\text{-}score = 0.838$. This was then used
747     to predict the taste of external molecules available in the FlavorDB, FooDB, SuperSweet (Ahmed *et*
748     *al.*, 2011), Super Natural II (Banerjee *et al.*, 2015), DSSTox (Richard & Williams, 2002) and
749     DrugBank libraries. In a further analysis, molecular taste of Bitter new, UNIMI set and Phytochemical
750     dictionary databases (Dagan-Wiener *et al.*, 2017) were also predicted. Fritz (Fritz *et al.*, 2021)
751     implemented the VirtualBitter model following the same workflow as for the sweetness prediction
752     (refer to the sweet and non-sweet section). They retrieved molecules from the BitterDB (Dagan-
753     Wiener *et al.*, 2019) and from their BitterSweetForest model (Banerjee & Preissner, 2018), in order

754 to model 1,612 bitterants and non-bitter (sweet and tasteless) molecules. The RF classifier exhibited

755 acceptable performances in prediction (20% test compounds): $NER = 0.898$, $AUC = 0.956$, $F\text{-score}$

756 $= 0.882$ and $ACC = 0.901$. In addition, when a molecule is predicted as bitterant, the webserver

757 provides the potential bitter target prediction for the 25 human bitter receptors (hTAS2Rs) based on

758 a similarity-based analysis. Finally, the VirtualBitter model was used to virtually screen diverse

759 molecules from the DrugBank database and the SuperNatural II database (Banerjee *et al.*, 2015).

760 Between 2020 and 2021, Charoenkwan's group published three webservers for taste prediction of a

761 curated database of 320 bitter peptides and 320 non-bitter peptides (BTP640), randomly generated

762 from BIOPEP (Minkiewicz *et al.*, 2008). The dataset was split into a training set and a test set (80:20).

763 NB and RF classifiers as well as several other classifiers were used: $k$NN, scoring card method

764 (SCM), bidirectional encoder representation from transformers (BERT), support vector machine

765 (SVM), decision tree (DT), extremely randomized trees (ETree), linear support vector classifier

766 (SVC), logistic regression (LR), multi-layer perceptron (MLP) and extreme gradient boosting (XGB).

767 The SCM classifier (Huang *et al.*, 2012), which was used through the dipeptide propensity score

768 (PDS) and optimized with GAs, achieved good results in prediction ($ACC = 0.844$, $AUC = 0.904$ and

769 $MCC = 0.688$) when compared to the SVM, RF, NB, $k$NN and DT, and it was included in the iBitter-

770 SCM webserver application (Charoenkwan *et al.*, 2020a). The authors stated that iBitter-SCM

771 constituted a useful tool for the high-throughput prediction and *de novo* design of novel bitterant

772 peptides. Another webserver, named BERT4Bitter (Charoenkwan *et al.*, 2021a), automatically

773 generates feature descriptors for peptides through the BERT algorithm. This model achieved the best

774 test set performance ($ACC = 0.922$, $AUC = 0.964$ and $MCC = 0.844$) with respect to the other

775 calibrated classifiers (DT, ETree, $k$NN, SVC, LR, MLP, NB, RF, SVM and XGB). For the webserver

776 iBitter-Fuse (Charoenkwan *et al.*, 2021b), five groups of molecular features were calculated: 20

777 amino acid composition (AAC), 400 dipeptide composition (DPC), 21 pseudo amino acid

778 composition (PAAC), 22 amphiphilic pseudo amino acid composition (APAAC), 531

779 physicochemical properties from AAindex (AAI), as well as a new group achieved by fusing features

780 (994 descriptors). Ten SVM models were calculated, providing excellent prediction quality ($ACC =$

781 $0.930$, $AUC = 0.933$ and $MCC = 0.859$). As described in their previous work, the authors calibrated

782 other machine learning models and demonstrated that the iBitter-Fuse model was superior in any case

783 (refer to Table 3 for the comparison between the iBitter-Fuse and the iBitter-SCM and BERT4Bitter

784 classifiers).

785 Dagan-Wiener used the Adaptive Boosting (AdaBoost) classifier for the first time in this framework

786 to create the BitterPredict model (Dagan-Wiener *et al.*, 2017). The dataset was composed of 691

787 bitterants (632 from the BitterDB (Wiener *et al.*, 2012)) and 1,917 non-bitter compounds retrieved

24

from several sources, which included 1,360 non-bitter flavors, 336 sweeteners, 186 tasteless molecules and 35 molecules labelled as non-bitter (molecules not described by the word bitter in the source). Each compound was represented by 59 molecular descriptors. The model was finally trained with 16 molecular descriptors and demonstrated predictive ability for the test set (30% of molecules) with $NER = 0.812$ and $ACC = 0.832$. Subsequently, the bitter class was evaluated for three external datasets, namely Bitter New ($Sn = 0.739$), UNIMI set ($Sn = 0.783$) and Phytochemical Dictionary (Baxter *et al.*, 1999) ($Sn = 0.980$ and $Sp = 0.692$). In a further step, the BitterPredict classifier was applied to achieve prospective predictions of compounds from the FooDB, DrugBank, ChEBI and the database of natural products. One year later, Zheng (Zheng *et al.*, 2018) developed several models based on the gradient boosting machine (GBM), as well as $k$NN, SVM, RF and two deep neuron networks (DNN2 and DNN3). These authors used a curated dataset of 707 bitterants and 592 non-bitter compounds (132 tasteless, 17 non-bitter and 443 sweet). Molecules were represented by means of several extended connectivity fingerprints (ECFPs): 1024bit-ECFP4, 2048bit-ECFP4, 1024bit-ECFP6 and 2048bit-ECFP6. In order to avoid bias due to partition, the splitting of the dataset was repeated 19 times for the $k$NN, SVM, GBM and RF models, and three times for the DNN2 and DNN3 models. Thus, 1,312 individual models and 96 average models were calibrated and consensus voting was used to obtain nine models (CM01 - CM09), which were integrated in the server e-Bitter tool. The best model (CM01) exhibited the following parameters for the test set (20% of compounds): *F-score* $= 0.936$, $ACC = 0.929$ and $MCC = 0.856$.

The XGBoost classifier was also used for bitterness prediction. Margulis proposed the BitterIntense model (Margulis *et al.*, 2021) for the classification of bitter molecules into very bitter and non-very bitter (including non-bitter) classes. A dataset of 721 compounds were obtained from behavioral studies using the rat brief access taste aversion (BATA), BitterDB (Dagan-Wiener *et al.*, 2019), Analyticon repository of natural compounds on Kaggle, as well as from their previous dataset BitterPredict (Dagan-Wiener *et al.*, 2017). Subsequently, 3D structures were used to calculate Canvas molecular descriptors (Schrödinger LLC, 2017) and QikProp features (ADME descriptors) (Schrödinger LLC, 2015). The XGBoost model achieved acceptable prediction on the test set (105 tastants): $NER = 0.790$, *F-score* $= 0.700$ and $ACC = 0.800$. Moreover, the BitterIntense model was used for analyzing the connection between toxicity and the level of bitterness of molecules, as well as for potential repurposing of COVID-19 targets. Independently, Bai developed the Children's Bitter Drug Prediction System' (CBDPS) (Bai *et al.*, 2021) for the bitterness prediction of medicines. The experimental dataset was retrieved from published works and the BitterDB (Dagan-Wiener *et al.*, 2019), which consisted of 1,732 tastants with a balanced number between bitter and non-bitter tastants (ratio of 1:1). Then, 166 MACCS structural keys and 114 ChemoPy descriptors (Cao *et al.*, 2013)

were used to calibrate four models based on the XGBoost and RF classifiers. Among these models, the optimal one was obtained with the XGBoost classifier and the MACCS structural keys, and achieved the following performance in cross-validation: *F-score* = 0.881 and *ACC* = 0.882. In a last step, the CBDPS model was applied to the screening of the external dataset of 222 children's oral medicines.

The XGBoost classifier was also applied to develop the BitterMatch model (Margulis *et al.*, 2022). A curated dataset of 303 bitterants resulted in 4,501 pairs of ligand-receptor associations (740 positives and 3,761 negatives). Optimized bitterants were used to calculate Canvas descriptors (Schrödinger LLC, 2017), while 3 sets of features were computed for receptors. The BitterMatch algorithm was divided into two scenarios: *filling the gaps* and *new ligands*. In both cases, 20% of the molecules were considered as test sets, keeping in mind the proportion of the classes (repeated 100 times). In *filling the gaps*, the best model included chemical properties and neighbor-informed chemical similarity features with an average recall-precision of 0.759. In contrast, the *new ligands* scenario considered only chemical properties of ligands and receptors, as well as neighbors-informed ligand similarity features (average recall-precision of 0.699). Afterwards, it was used to predict associations for 12 external bitterants and drugs from the DrugBank.

More recently, Bo (Bo *et al.*, 2022) calibrated quantitative structure-taste relationships based on MLP and CNN deep learning classifiers (as described before in the sweet/bitter and sweet/non-sweet sections). In this work, a dataset of 797 bitterants and 1,436 non-bitterants (sweet and tasteless) was used. The BitterMLP-Descriptor classifier with seven RDKit descriptors exhibited similar validation performance (*NER* = 0.820, *AUC* = 0.940, *ACC* = 0.840 and *MCC* = 0.660) with respect to the BitterCNN classifier (*NER* = 0.790, *AUC* = 0.880, *ACC* = 0.810 and *MCC* = 0.600). As described in the sweet/non-sweet models, these two classifiers were used to analyze 1,229 tastants from the sweet data set. De León (De León *et al.*, 2022) calibrated SVM, RF, AdaBoost and *k*NN models for a curated dataset of 932 bitterants and 1,908 non-bitter molecules retrieved from BitterDB (Dagan-Wiener *et al.*, 2019), Fenaroli's Handbook of flavours (Burdock, 2010) and the dataset of Rojas (Rojas *et al.*, 2016c). The compounds were represented by ECFPs and 22 selected Mordred descriptors (Moriwaki *et al.*, 2018) on the basis of their probability density. For validation purposes, 20% of the molecules were included in the test set. The two best classifiers turned out to be SVM ($ACC_{train}$ = 0.836 and $ACC_{test}$ = 0.870) and AdaBoost ($ACC_{train}$ = 0.842 and $ACC_{test}$ = 0.847) based on ECFPs and descriptors, respectively. In addition, the UNIMI dataset (Dagan-Wiener *et al.*, 2017) was used as the external set to validate the performance of Premexotac models.

**2.4. Umami and non-umami tastants**

There are fewer ligand-based (LB) machine learning models that have been developed for the discrimination between umami and non-umami peptides. This could be due to the higher complexity of sensory evaluation and related costs than those related to the evaluation of sweet and bitter molecules.

For the first model, named iUmami-SCM (Charoenkwan *et al.*, 2020b), the experimental information for umami peptides was retrieved from the literature and from the BIOPEP-UWM database, while bitter peptides, previously studied by the authors, were considered as non-umami molecules. The UMP442 database (140 umami and 302 non-umami peptides) was used to calibrate a SCM classifier based on a dipeptide propensity score (PDS), as described in the iBitter-SCM model (Charoenkwan *et al.*, 2020a). The best model achieved good results in prediction (20% of test molecules): $AUC = 0.898$, $ACC = 0.865$, $MCC = 0.679$, $Sn = 0.714$ and $Sp = 0.934$. In addition, the model's performance was compared with six ML classifiers (SVM, RF, MLP, NB, $k$NN and DT). In the second application, the same group of Charoenkwan combined six well-known ML classifiers (ETree, $k$NN, LR, PLS, RF and SVM) in the UMPred-FRL model (Charoenkwan *et al.*, 2021c). To this end, they used molecules of the UMP442 database (Charoenkwan *et al.*, 2020b), which were represented by seven feature descriptors: amino acid composition (AAC), amphiphilic pseudo-amino acid composition (APAAC), dipeptide composition (DPC), composition (CTDC), transition (CTDT), distribution (CTDD) and pseudo-amino acid composition (PAAC). The UMPred-FRL predictor was assembled by the best 7 informative features (SVM-AAC, PLS-AAC, SVM-CTDC, RF-DPC, RF-CTDC, PLS-APAAC and LRDPC), and exhibited better performances when compared to the iUmami-SCM classifier prediction ($AUC = 0.919$, $ACC = 0.888$ and $MCC = 0.735$).

In 2022, Pallante developed the Virtuous Umami platform (Pallante *et al.*, 2022) for umami prediction based on SVM classifiers and the Charoenkwan's UMP442 database (Charoenkwan *et al.*, 2020b). Due to the unbalanced classes, umami peptides were randomly duplicated to balance the class cardinalities. Subsequently, 1,613 conformation-independent Mordred features (Moriwaki *et al.*, 2018) were subjected to feature selection by means of different approaches, which were used to calibrate diverse SVM models. The best prediction was achieved by consensus between two models (12 features), which exhibited a slightly lower performance in prediction ($AUC = 0.850$, $F\text{-}score = 0.793$ and $ACC = 0.876$) when compared to the iUmami-SCM and UMPred-FRL predictors. The effectiveness of the model was visually shown by means of t-SNE. Finally, the umami predictor was used to virtually screen the FooDB, FlavorDB, PhenolExplorer, Natural Product Atlas and PhytoHub databases.

Recently, Dutta proposed the identification of optimal sequential residue patterns for umami and bitter peptides (Dutta *et al.*, 2022a). These authors used a curated database of 292 bitter and 146

umami compounds retrieved from Charoenkwan's UMP442 database (Charoenkwan *et al.*, 2020b) and others sources. Each peptide was represented by the following coarse-grained representation: hydrophobic (H), polar and hydrophilic (P), positively charged (+) and negatively charged (−). Afterwards, seven libraries of peptides were created by repeating a fixed set of coarse-grained patterns. To select the best length, the dataset of taste-labeled peptides was split into a training set (80%) and a test set (20%) by means of stratified random sampling. A length of five ($N = 5$) was selected as the best coarse-grained pattern, where bitter peptides were represented by one hydrophobic followed by four polar residues (HPPPP), while umami peptides had two negative followed by three polar residues (−−PPP). In a further step, the authors tested this method by using two bitter proteins (Patatin-T5 and Legumin-A), where 8 and 5 peptide sequences with the aforementioned course-grain pattern were identified. This approach allowed the rapid screening and identification of sequential information patterns hidden in long chain peptides and proteins, rather than predicting the taste class of peptides (no classification performances were reported).

**2.5 Bitter, sweet and umami tastant**

Dutta developed the first deep learning classifier to discriminate among sweet, bitter and umami tastants (Dutta *et al.*, 2022b). The curated dataset was composed of 1,938 bitterants, 2,079 sweeteners and 98 umami compounds, which were retrieved from the ChemTastesDB database (Rojas *et al.*, 2022) and the BitterSweet dataset (Tuwani *et al.*, 2019). Afterwards, 102 RDKit molecular descriptors were used after a filtering process. For pattern recognition, the authors developed two chemical spaces based on PCA and t-SNE, along with a functional group analysis by computing the frequency of predefined fragments. Then, a deep neural network (DNN) with two hidden layers of 100 neurons was trained with 200 epochs. For balancing the cardinality of the umami class, the synthetic minority oversampling technique (SMOTE) for data augmentation was used. The DNN model was interpreted by means of the Shapley additive explanations (SHAP). The DNN model achieved good predictive performance (15% of compounds): $NER = 0.901$ and $ACC = 0.887$. Moreover, a graph neural network (GNN) was also tested with a slightly lower quality on external prediction ($NER = 0.865$ and $ACC = 0.896$). Independently, Xiu (Xiu *et al.*, 2022) used the BitterSweet dataset (Tuwani *et al.*, 2019) to develop the PyUmami model, which combined sweet and bitter classifiers based on multilayer perception (MLP) and Mordred descriptors. Then, the sweet-MLP ($ACC = 0.830$ and $AUC = 0.897$) and bitter-MLP models ($ACC = 0.81$ and $AUC = 0.895$) were used to predict the sweetness of 1,040 bitterants from the BitterDB, and the bitterness of 14,175 sweeteners from the SWEET-DB, respectively. Only 169 tastants predicted as both sweet/bitter by the PyUmami model were submitted to docking analysis with the T1R2/T1R3 and hT2R1 receptors.

922     Finally, 18 targets were experimentally verified for sweet, bitter and umami intensities by means of

923     electronic tongue analysis, and only 8 tastants were predicted to be non-toxic by means of twelve

924     QSAR approaches and three virtual Adverse Outcome Pathway (vAOP) models.

**2.6. Sour and non-sour tastants**

926     Only one LB classifier for the discrimination between sour and non-sour compounds has been

927     proposed (Fritz *et al.*, 2021). Information of molecules was retrieved from ChEMBL (Gaulton *et al.*,

928     2012) and curated from the PubMed database (Sayers *et al.*, 2021). The dataset consisted of 1,347

929     compounds divided into a training set and a test set of 1,214 and 133 molecules, respectively. The

930     model, named VirtualSour, was a ligand-based approach considering the RF classifier integrated with

931     the augmented random data sampling method. The model achieved good results in cross-validation

932     ($NER = 0.955$, $AUC = 0.998$, $F\text{-}score = 0.980$ and $ACC = 0.978$,) and prediction ($NER = 0.896$, $AUC$

933     $= 0.994$, $F\text{-}score = 0.842$ and $ACC = 0.977$).

**2.7 General trends in taste modelling**

935     When looking at the evolution of modelling approaches for predicting the different tastes, common

936     trends and tendencies can be seen. Figure 3 shows the number of molecules (included in both training

937     and test set) used for the development of structure-property models as a function of the publication

938     year, starting from the very beginning of the modelling era (1980) up to 2022. First of all, it is apparent

939     that the number of chemicals used to train or test QSAR models has greatly increased (note that the

940     y axis of Figure 3 is in log10 units). While the first modeling attempts considered a few dozen

941     chemicals, the number increased to several hundreds from 2000 to 2010. In addition, models were

942     initially developed considering only small families of compounds (for instance aldoximes,

943     perillartines, aspartyl dipeptides and sulfamates), which established restricted chemical spaces for

944     only these types of compounds. The most relevant increase in the number of chemicals occurred after

945     2015, when scientists started to use several thousand molecules to develop new models for taste

946     prediction. Interestingly, the study for bitter prediction published by Rodgers in 2006 (Rodgers *et al.*,

947     2006) used 13,530 molecules randomly selected from the MDL Drug Data Repository under the

948     assumption that this was representative of the bitterness chemical space. However, these molecules

949     were not validated with experimental sensory data as considered in the most recently published works.

**Figure 3 should be inserted around here**

The large growth of the number of molecules used in the development of models that started in 2015 is probably due to several research groups who concentrated their efforts on creating more extensive and comprehensive databases, such as SuperSweet (Ahmed *et al.*, 2011), BitterDB (Wiener *et al.*, 2012), ChEMBL (Gaulton *et al.*, 2012) and Super Natural II (Banerjee *et al.*, 2015). These databases collected and cataloged a greater number of substances associated with their molecular structures and experimental taste values, which enabled the subsequent development of models based on a significantly higher number of chemicals in the years after 2015. These large databases included heterogeneous molecules, which allowed the extension of chemical spaces and, in fact, some attempts were made for virtual screening of potential new tastants in several available databases, which were complemented, in some cases, with docking analysis and experimental sensory evaluation of the elicited tastants.

Another general trend is related to the type of analyzed taste. Figure 3, shows that in the first 29 years from the first model developed in 1980, sweetness was the principal interest. Within this modelling framework, only models for the discrimination of sweet chemicals versus bitter or non-sweet molecules were taken into account. Afterwards, due to the development of more comprehensive databases, models for the prediction of bitterness were proposed in addition to sweetness. The interest in bitterness prediction could be related to the increasing interest of using bitterants as food and pharmaceutical additives along with other applications. Starting from 2020, umami prediction proved to be another attractive topic in the scientific community. The increasing interest in modelling this taste is mainly related to Asian research groups, due to the importance of umami in oriental gastronomy. On the other hand, modelling of sourness and saltiness is limited by the reduced number of molecules that imprint these tastes.

The increasing number of molecules used to model tastants also enabled a better estimation of predictive performance; that is, the accuracy in the prediction of the taste of chemicals which were not used for model training. Validation is fundamental in the development of QSARs and usually consists of the use of some chemicals, with known experimental taste values but not involved in the model training, as the test molecules. The first studies did not generally account for model validation. Until 2016 less than 10 chemicals were used in a couple of studies to validate models for discrimination between sweet and bitter tastants (Table 1), while for the classification of sweet and non-sweet chemicals, no test compounds were considered until 2005 and just a few in the studies published between 2006 and 2009 (Table 2). On the other hand, the number of substances used for model validation has grown enormously in recent years and now, hundreds of molecules are normally used for validation purposes.

Finally, the increasing availability of newly synthesized chemicals has influenced the type of machine learning approaches that have been used to establish molecular structure-taste relationships. Initially only simple classification algorithms were used (such as Discriminant Analysis and CART), whereas in the last decade, advanced approaches have been frequently applied, such as RF, SVM, boosting algorithms and Neural Networks. This is a general trend in the framework of machine learning, which has been supported by the computational and technological advancements of the latest decades. However, unlike traditional approaches, the newer and novel classification methods require a tuning phase for the selection of optimal values of their hyperparameters. This tuning phase is executed by optimizing the models on a further set of chemicals, usually named an evaluation set, which has to be added to the training set (used for the learning phase) and the test set (used for the final validation phase). Therefore, execution of the tuning phase requires a more extended number of chemicals for their calculation.

It is interesting to note that although a very limited number of descriptors was used in the first developed models, the evolution of modeling approaches has not caused a considerable increase in the complexity of the models. In many cases the total number of descriptors used for the development of models is measured in the 10s, and only a few hundred descriptors have been used in some models for the discrimination of sweet and non-sweet tastants. Of course, molecular fingerprints are a special case, since the thousands of binary bits they include have to be considered simultaneously as a holistic description of the molecular structure. As in other modelling frameworks, the limited number of descriptors is probably due to the maintenance of a correct balance between the model complexity, predictive ability and interpretability.

In earlier models for sweet prediction, descriptors mainly related to molecular size and bulkiness were used while recently, quantum-chemical descriptors were considered as well as different types of fingerprints and descriptors calculated by means of different software including Dragon, RDKit, Mordred, Pybel, alvaDesc and MOE2d. From an analysis of the most frequent molecular descriptors in models for bitterness, the relevant structural features are the presence of carbon/oxygen groups, sugar moieties, quaternary carbon centers and highly branched carbon centers, physicochemical properties, specific properties of the molecular surface and hydrophobicity. More specifically, the bitterness of peptides is strongly related to composition of amino acids, dipeptides and pseudo amino acids. Finally, molecular descriptors used for modelling umami taste are mainly linked to the presence of hydrophilic amino acids with negative charge and low molecular weights. In addition, patterns in the scaffolds related to amino acid composition; specifically glutamic acid (Glu) and aspartic acid (Asp) amino acid, were found to be crucial for umami prediction of peptides.

## 3. Conclusions

In this paper, we present a logical, comprehensive and critical review of the current state of ligand-based models of quantitative structure-property relationships along with the history of the prediction of the taste of molecules. Models detailed here complement previously published reviews available in the literature. Although the main modeling applications presented in this review relate to the prediction of molecular sweetness and bitterness, there is a notable increase in the interest and proposed application of QSAR models for the prediction of umami and sour tastants. It is notable that many authors cited in this review attempted to use the largest possible databases of tastants, as well as to improve the chemical representation of these databases through the use of several molecular descriptors, structural keys and fingerprints. In addition, this review reflects the wide variety of machine learning approaches used by investigators in order to calibrate more general models used in the prediction of properties of new molecules. In the future, it is expected that *in silico* methods will increase the application of predictive models in food chemistry (foodinformatics) in order to better understand the mechanisms involved in taste prediction. In addition, predictive models may provide useful tools to discover new molecular tastants with potential uses as raw-materials or additives in the food and pharmaceutical industries. Finally, our recommendation to chemists involve in taste prediction is to develop the largest possible molecular tastant databases to be used with novel classifiers in order to develop models able to predict more than two classes at a time. This expanded capability will greatly advance the science of foodinformatics.

## 4. Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## 5. Acknowledgements

## 6. References

Acton, E. M., & Stone, H. (1976). Potential new artificial sweetener from study of structure-taste relationships. *Science, 193*(4253), 584-586. https://doi.org/10.1126/science.959816

1048     Adler, E., Hoon, M. A., Mueller, K. L., Chandrashekar, J., Ryba, N. J., & Zuker, C. S. (2000). A
1049          novel family of mammalian taste receptors. *Cell, 100*(6), 693-702.
1050          https://doi.org/10.1016/S0092-8674(00)80705-9

1051     Ahmed, J., Preissner, S., Dunkel, M., Worth, C. L., Eckert, A., & Preissner, R. (2011). SuperSweet-
1052          A resource on natural and artificial sweetening agents. *Nucleic Acids Research, 39*, D377-
1053          D382. https://doi.org/10.1093/nar/gkq917

1054     Bai, G., Wu, T., Zhao, L., Wang, X., Li, S., & Ni, X. (2021). CBDPS 1.0: A Python GUI application
1055          for machine learning models to predict bitter-tasting children's oral medicines. *Chemical and*
1056          *Pharmaceutical Bulletin, 69*, 989-994. https://doi.org/10.1248/cpb.c20-00866

1057     Baines, D., & Brown, M. (2016). Flavor enhancers: Characteristics and uses. In B. Caballero, P. M.
1058          Finglas & F. Toldrá (Eds.), *Encyclopedia of food and health*, vol. 2 (pp. 716-723). Oxford
1059          (UK): Academic Press.

1060     Ballabio, D., Consonni, V., Mauri, A., Claeys-Bruno, M., Sergent, M., & Todeschini, R. (2014). A
1061          novel variable reduction method adapted from space-filling designs. *Chemometrics and*
1062          *Intelligent Laboratory Systems, 136*, 147-154.
1063          https://doi.org/10.1016/j.chemolab.2014.05.010

1064     Ballabio, D., Grisoni, F., & Todeschini, R. (2018). Multivariate comparison of classification
1065          performance measures. *Chemometrics and Intelligent Laboratory Systems, 174*, 33-44.
1066          https://doi.org/10.1016/j.chemolab.2017.12.004

1067     Banerjee, P., Erehman, J., Gohlke, B.-O., Wilhelm, T., Preissner, R., & Dunkel, M. (2015). Super
1068          Natural II-A database of natural products. *Nucleic Acids Research, 43*(D1), D935-D939.
1069          https://doi.org/10.1093/nar/gku886

1070     Banerjee, P., & Preissner, R. (2018). BitterSweetForest: A random forest based binary classifier to
1071          predict bitterness and sweetness of chemical compounds. *Frontiers in Chemistry, 6*, 93.
1072          https://doi.org/10.3389/fchem.2018.00093

1073     Barker, M., & Rayens, W. (2003). Partial least squares for discrimination. *Journal of Chemometrics,*
1074          *17*, 166-173. https://doi.org/10.1002/cem.785

1075     Baxter, H., Harborne, J. B., & Moss, G. P. (1999). *Phytochemical dictionary: A handbook of bioactive*
1076          *compounds from plants* (Second ed.). Boca Raton, USA: CRC press.

1077     Bayer, S., Mayer, A. I., Borgonovo, G., Morini, G., Di Pizio, A., & Bassoli, A. (2021).
1078          Chemoinformatics view on bitter taste receptor agonists in food. *Journal of Agricultural and*
1079          *Food Chemistry, 69*(46), 13916-13924. https://doi.org/10.1021/acs.jafc.1c05057

1080     Behrens, M., & Ziegler, F. (2020). Structure-function analyses of human bitter taste receptors-where
1081          do we stand? *Molecules, 25*(19), 4423. https://doi.org/10.3390/molecules25194423

1082 Bender, A., Mussa, H. Y., Glen, R. C., & Reiling, S. (2004). Similarity searching of chemical
1083     databases using atom environment descriptors (MOLPRINT 2D): evaluation of performance.
1084     *Journal of Chemical Information and Computer Sciences, 44*(5), 1708-1718.
1085     https://doi.org/10.1021/ci0498719

1086 Bo, W., Qin, D., Zheng, X., Wang, Y., Ding, B., Li, Y., & Liang, G. (2022). Prediction of bitterant
1087     and sweetener using structure-taste relationship models based on an artificial neural network.
1088     *Food Research International, 153*, 110974. https://doi.org/10.1016/j.foodres.2022.110974

1089 Breiman, L. (2001). Random forests. *Machine learning, 45*, 5-32.

1090 Breiman, L. J., Friedman, J. H., Olsen, R., & Stone, C. (1984). *Classification and regression trees*.
1091     California (USA): Wadsworth, Belmont

1092 Brereton, R. G. (2011). One- class classifiers. *Journal of Chemometrics, 25*, 225-246.
1093     https://doi.org/10.1002/cem.1397

1094 Brereton, R. G., & Lloyd, G. R. (2014). Partial least squares discriminant analysis: Taking the magic
1095     away. *Journal of Chemometrics, 28*, 213-225. https://doi.org/10.1002/cem.2609

1096 Breslin, P., & Huang, L. (2006). Human Taste: Peripheral Anatomy, tastetransduction, and coding.
1097     In T. Hummel & A. Welge-Lüssen (Eds.), *Taste and smell*, (pp. 152-190): KARGER.

1098 Brockhoff, A., Behrens, M., Roudnitzky, N., Appendino, G., Avonto, C., & Meyerhof, W. (2011).
1099     Receptor agonism and antagonism of dietary bitter compounds. *Journal of Neuroscience,*
1100     *31*(41), 14775-14782. https://doi.org/10.1523/JNEUROSCI.2923-11.2011

1101 Burdock, G. A. (2010). *Fenaroli's handbook of flavor ingredients*. Boca Ratón, USA: CRC Press.

1102 Cao, D.-S., Xu, Q.-S., Hu, Q.-N., & Liang, Y.-Z. (2013). ChemoPy: Freely available python package
1103     for computational biology and chemoinformatics. *Bioinformatics, 29*(8), 1092-1094.
1104     https://doi.org/10.1093/bioinformatics/btt105

1105 Chandrashekar, J., Hoon, M. A., Ryba, N. J. P., & Zuker, C. S. (2006). The receptors and cells for
1106     mammalian taste. *Nature, 444*, 288-294. https://doi.org/10.1038/nature05401

1107 Charoenkwan, P., Yana, J., Schaduangrat, N., Nantasenamat, C., Hasan, M. M., & Shoombuatong,
1108     W. (2020a). iBitter-SCM: Identification and characterization of bitter peptides using a scoring
1109     card method with propensity scores of dipeptides. *Genomics, 112*(4), 2813-2822.
1110     https://doi.org/10.1016/j.ygeno.2020.03.019

1111 Charoenkwan, P., Yana, J., Nantasenamat, C., Hasan, M. M., & Shoombuatong, W. (2020b).
1112     iUmami-SCM: A novel sequence-based predictor for prediction and analysis of umami
1113     peptides using a scoring card method with propensity scores of dipeptides. *Journal of*
1114     *Chemical Information and Modeling, 60*(12), 6666-6678.
1115     https://doi.org/10.1021/acs.jcim.0c00707

Charoenkwan, P., Nantasenamat, C., Hasan, M. M., Manavalan, B., & Shoombuatong, W. (2021a). BERT4Bitter: A bidirectional encoder representations from transformers (BERT)-based model for improving the prediction of bitter peptides. *Bioinformatics, 37*(17), 2556-2562. https://doi.org/10.1093/bioinformatics/btab133

Charoenkwan, P., Nantasenamat, C., Hasan, M., Moni, M. A., Lio, P., & Shoombuatong, W. (2021b). iBitter-fuse: A novel sequence-based bitter peptide predictor by fusing multi-view features. *International Journal of Molecular Sciences, 22*(16), 8958. https://doi.org/10.3390/ijms22168958

Charoenkwan, P., Nantasenamat, C., Hasan, M. M., Moni, M. A., Manavalan, B., & Shoombuatong, W. (2021c). UMPred-FRL: A new approach for accurate prediction of umami peptides using feature representation learning. *International Journal of Molecular Sciences, 22*, 13124. https://doi.org/10.3390/ijms222313124

Chattopadhyay, S., Raychaudhuri, U., & Chakraborty, R. (2014). Artificial sweeteners–A review. *Journal of Food Science and Technology, 51*, 611-621. https://doi.org/10.1007/s13197-011-0571-1

Chaudhari, N., Pereira, E., & Roper, S. D. (2009). Taste receptors for umami: The case for multiple receptors. *The American Journal of Clinical Nutrition, 90*(3), 738S-742S. https://doi.org/10.3945/ajcn.2009.27462H

Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. *KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data*, 785-794. https://doi.org/10.1145/2939672.2939785

Chéron, J.-B., Casciuc, I., Golebiowski, J., Antonczak, S., & Fiorucci, S. (2017). Sweetness prediction of natural compounds. *Food Chemistry, 221*, 1421-1425. https://doi.org/10.1016/j.foodchem.2016.10.145

Dagan-Wiener, A., Nissim, I., Abu, N. B., Borgonovo, G., Bassoli, A., & Niv, M. Y. (2017). Bitter or not? BitterPredict, a tool for predicting taste from chemical structure. *Scientific Reports, 7*, 12074. https://doi.org/10.1038/s41598-017-12359-7

Dagan-Wiener, A., Di Pizio, A., Nissim, I., Bahia, M. S., Dubovski, N., Margulis, E., & Niv, M. Y. (2019). BitterDB: Taste ligands and receptors database in 2019. *Nucleic Acids Research, 47*(D1), D1179-D1185. https://doi.org/10.1093/nar/gky974

Damodaran, S., & Parkin, K. L. (2017). *Fennema's food chemistry* (5th ed.). Boca Raton (USA): CRC Press.

De León, G., Fröhlich, E., & Salar-Behzadi, S. (2021). Bitter taste *in silico*: A review on virtual ligand screening and characterization methods for TAS2R-bitterant interactions. *International Journal of Pharmaceutics, 600*, 120486. https://doi.org/10.1016/j.ijpharm.2021.120486

De León, G., Fröhlich, E., Fink, E., Di Pizio, A., & Salar-Behzadi, S. (2022). Premexotac: Machine learning bitterants predictor for advancing pharmaceutical development. *International Journal of Pharmaceutics, 628*, 122263. https://doi.org/10.1016/j.ijpharm.2022.122263

Deng, X., Lin, H., Ahmed, I., & Sui, J. (2021). Isolation and identification of the umami peptides from Trachinotus ovatus hydrolysate by consecutive chromatography and Nano-HPLC-MS/MS. *LWT-Food Science and Technology, 141*, 110887. https://doi.org/10.1016/j.lwt.2021.110887

DeSimone, J. A., & Lyall, V. (2006). Taste receptors in the gastrointestinal tract III. Salty and sour taste: sensing of sodium and protons by the tongue. *American Journal of Physiology-Gastrointestinal and Liver Physiology, 291*(6), G1005-G1010. https://doi.org/10.1152/ajpgi.00235.2006

Di Lorenzo, P. M., Chen, J.-Y., Rosen, A. M., & Roussin, A. T. (2009). Tastant. In M. D. Binder, N. Hirokawa & U. Windhorst (Eds.), *Encyclopedia of neuroscience*, (pp. 4014-4019). Berlin (Germany): Springer.

Di Pizio, A., & Niv, M. Y. (2015). Promiscuity and selectivity of bitter molecules and their receptors. *Bioorganic & Medicinal Chemistry, 23*(14), 4082-4091. https://doi.org/10.1016/j.bmc.2015.04.025

Doty, R. L., Bagla, R., Morgenson, M., & Mirza, N. (2001). NaCl thresholds: Relationship to anterior tongue locus, area of stimulation, and number of fungiform papillae. *Physiology & Behavior, 72*(3), 373-378. https://doi.org/10.1016/S0031-9384(00)00416-9

Drew, M. G. B., Wilden, G. R. H., Spillane, W. J., Walsh, R. M., Ryder, C. A., & Simmie, J. M. (1998). Quantitative structure-activity relationship studies of sulfamates $RNHSO_3Na$: Distinction between sweet, sweet-bitter, and bitter molecules. *Journal of Agricultural and Food Chemistry, 46*(8), 3016-3026. https://doi.org/10.1021/jf980095c

Durant, J. L., Leland, B. A., Henry, D. R., & Nourse, J. G. (2002). Reoptimization of MDL keys for use in drug discovery. *Journal of Chemical Information and Computer Sciences, 42*(6), 1273-1280. https://doi.org/10.1021/ci010132r

Dutta, A., Bereau, T., & Vilgis, T. A. (2022a). Identifying sequential residue patterns in bitter and umami peptides. *ACS Food Science & Technology, 2*(11), 1773-1780. https://doi.org/10.1021/acsfoodscitech.2c00251

Dutta, P., Jain, D., Gupta, R., & Rai, B. (2022b). Classification of tastants: A deep learning based approach. *ChemRxiv*. https://doi.org/10.26434/chemrxiv-2022-rs6x3

Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences, 55*(1), 119-139. https://doi.org/10.1006/jcss.1997.1504

Fritz, F., Preissner, R., & Banerjee, P. (2021). VirtualTaste: A web server for the prediction of organoleptic properties of chemical compounds. *Nucleic Acids Research, 49*(W1), W679-W684. https://doi.org/10.1093/nar/gkab292

Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., & Al-Lazikani, B. (2012). ChEMBL: A large-scale bioactivity database for drug discovery. *Nucleic Acids Research, 40*(D1), D1100-D1107. https://doi.org/10.1093/nar/gkr777

Gilbertson, T. A., Fontenot, D. T., Liu, L., Zhang, H., & Monroe, W. T. (1997). Fatty acid modulation of $K^+$ channels in taste receptor cells: Gustatory cues for dietary fat. *American Journal of Physiology-Cell Physiology, 272*(4), C1203-C1210. https://doi.org/10.1152/ajpcell.1997.272.4.C1203

Goel, A., Gajula, K., Gupta, R., & Rai, B. (2021). In-silico screening of database for finding potential sweet molecules: A combined data and structure based modeling approach. *Food Chemistry, 343*, 128538. https://doi.org/10.1016/j.foodchem.2020.128538

Gramatica, P. (2007). Principles of QSAR models validation: Internal and external. *QSAR & Combinatorial Science, 26*(5), 694-701. https://doi.org/10.1002/qsar.200610151

Gu, J., Gui, Y., Chen, L., Yuan, G., Lu, H.-Z., & Xu, X. (2013). Use of natural products as chemical library for drug discovery and network pharmacology. *PloS one, 8*(4), e62839. https://doi.org/10.1371/journal.pone.0062839

Hand, D. J. (1997). *Construction and assessment of classification rules*. Chichester (UK): Wiley.

Huang, H.-L., Charoenkwan, P., Kao, T.-F., Lee, H.-C., Chang, F.-L., Huang, W.-L., Ho, S.-J., Shu, L.-S., Chen, W.-L., & Ho, S.-Y. (2012). Prediction and analysis of protein solubility using a novel scoring card method with dipeptide composition. *BMC Bioinformatics, 13*(S17), S3. https://doi.org/10.1186/1471-2105-13-S17-S3

Huang, W., Shen, Q., Su, X., Ji, M., Liu, X., Chen, Y., Lu, S., Zhuang, H., & Zhang, J. (2016). BitterX: A tool for understanding bitter taste in humans. *Scientific Reports, 6*, 23450. https://doi.org/10.1038/srep23450

Kelly, D. P., Spillane, W. J., & Newell, J. (2005). Development of structure-taste relationships for monosubstituted phenylsulfamate sweeteners using classification and regression tree (CART)

analysis. *Journal of Agricultural and Food Chemistry, 53*(17), 6750-6758. https://doi.org/10.1021/jf0507137

Kier, L. B. (1980). Molecular structure influencing either a sweet or bitter taste among aldoximes. *Journal of Pharmaceutical Sciences, 69*(4), 416-419. https://doi.org/10.1002/jps.2600690414

Kode srl. (2018). Dragon version 7. Software for molecular descriptor calculation, http://chm.kode-solutions.net/

Kowalski, B., & Bender, C. (1972). *k*-Nearest Neighbor classification rule (pattern recognition) applied to nuclear magnetic resonance spectral interpretation. *Analytical Chemistry, 44*(8), 1405-1411. https://doi.org/10.1021/ac60316a008

Lavine, B. K., & Rayens, W. S. (2009). 3.27 - Classification: Basic Concepts. In S. Brown, R. Tauler & B. Walczak (Eds.), *Comprehensive Chemometrics (Second Edition)*, (pp. 567-573). Oxford, England: Elsevier. https://doi.org/10.1016/B978-0-444-64165-6.01010-7

Lee, J., Song, S. B., Chung, Y. K., Jang, J. H., & Huh, J. (2022). BoostSweet: Learning molecular perceptual representations of sweeteners. *Food Chemistry, 383*, 132435. https://doi.org/10.1016/j.foodchem.2022.132435

Ley, J., Reichelt, K., Obst, K., Krammer, G., & Engel, K. H. (2012). Important tastants and new developments. In H. Jeleń (Ed.), *Food flavors: Chemical, sensory and technological properties*, (pp. 19-33). Boca Raton (USA): CRC Press.

Liang, L., Duan, W., Zhang, J., Huang, Y., Zhang, Y., & Sun, B. (2022a). Characterization and molecular docking study of taste peptides from chicken soup by sensory analysis combined with nano-LC-Q-TOF-MS/MS. *Food Chemistry, 383*, 132455. https://doi.org/10.1016/j.foodchem.2022.132455

Liang, L., Zhou, C., Zhang, J., Huang, Y., Zhao, J., Sun, B., & Zhang, Y. (2022b). Characteristics of umami peptides identified from porcine bone soup and molecular docking to the taste receptor T1R1/T1R3. *Food Chemistry, 387*, 132870. https://doi.org/10.1016/j.foodchem.2022.132870

Liu, Z., Zhu, Y., Wang, W., Zhou, X., Chen, G., & Liu, Y. (2020). Seven novel umami peptides from Takifugu rubripes and their taste characteristics. *Food Chemistry, 330*, 127204. https://doi.org/10.1016/j.foodchem.2020.127204

Malavolta, M., Pallante, L., Mavkov, B., Stojceski, F., Grasso, G., Korfiati, A., Mavroudi, S., Kalogeras, A., Alexakos, C., & Martos, V. (2022). A survey on computational taste predictors. *European Food Research and Technology, 248*, 2215-2235. https://doi.org/10.1007/s00217-022-04044-5

1247 Margulis, E., Dagan-Wiener, A., Ives, R. S., Jaffari, S., Siems, K., & Niv, M. Y. (2021). Intense
1248     bitterness of molecules: Machine learning for expediting drug discovery. *Computational and*
1249     *Structural Biotechnology Journal, 19*, 568-576. https://doi.org/10.1016/j.csbj.2020.12.030

1250 Margulis, E., Slavutsky, Y., Lang, T., Behrens, M., Benjamini, Y., & Niv, M. Y. (2022). BitterMatch:
1251     Recommendation systems for matching molecules with bitter taste receptors. *Journal of*
1252     *Cheminformatics, 14*(1), 45. https://doi.org/10.1186/s13321-022-00612-9

1253 Maroni, G., Pallante, L., Di Benedetto, G., Deriu, M. A., Piga, D., & Grasso, G. (2022). Informed
1254     classification of sweeteners/bitterants compounds via explainable machine learning. *Current*
1255     *Research in Food Science, 5*, 2270-2280. https://doi.org/10.1016/j.crfs.2022.11.014

1256 Martinez-Mayorga, K., & Medina-Franco, J. L. (2014). *Foodinformatics: Applications of chemical*
1257     *information to food chemistry*. Cham (Switzerland): Springer.

1258 Mathea, M., Klingspohn, W., & Baumann, K. (2016). Chemoinformatic classification methods and
1259     their applicability domain. *Molecular Informatics, 35*, 160-180.
1260     https://doi.org/10.1002/minf.201501019

1261 Matsunami, H., Montmayeur, J.-P., & Buck, L. B. (2000). A family of candidate taste receptors in
1262     human and mouse. *Nature, 404*, 601-604. https://doi.org/10.1038/35007072

1263 Mauri, A., & Bertola, M. (2022). Alvascience: A New Software Suite for the QSAR Workflow
1264     Applied to the Blood-Brain Barrier Permeability. *International Journal of Molecular*
1265     *Sciences, 23*(21), 12882. https://doi.org/10.3390/ijms232112882

1266 McLachlan, G. J. (1992). *Discriminant analysis and statistical pattern recognition*. New York
1267     (USA): Wiley.

1268 Minkiewicz, P., Dziuba, J., Iwaniak, A., Dziuba, M., & Darewicz, M. (2008). BIOPEP database and
1269     other programs for processing bioactive peptide sequences. *Journal of AOAC International,*
1270     *91*(4), 965-980. https://doi.org/10.1093/jaoac/91.4.965

1271 Miyashita, Y., Takahashi, Y., Takayama, C., Sumi, K., Nakatsuka, K., Ohkubo, T., Abe, H., & Sasaki,
1272     S.-i. (1986a). Structure-taste correlation of L-Aspartyl dipeptides using the SIMCA method.
1273     *Journal of Medicinal Chemistry, 29*(6), 906-912. https://doi.org/10.1021/jm00156a006

1274 Miyashita, Y., Takahashi, Y., Takayama, C., Ohkubo, T., Funatsu, K., & Sasaki, S.-i. (1986b).
1275     Computer-assisted structure/taste studies on sulfamates by pattern recognition methods.
1276     *Analytica Chimica Acta, 184*, 143-149. https://doi.org/10.1016/S0003-2670(00)86477-6

1277 Morgan, H. L. (1965). The generation of a unique machine description for chemical structures-a
1278     technique developed at chemical abstracts service. *Journal of Chemical Documentation, 5*(2),
1279     107-113. https://doi.org/10.1021/c160017a018

Morini, G., Bassoli, A., & Borgonovo, G. (2011). Molecular modelling and models in the study of sweet and umami taste receptors. A review. *Flavour and Fragrance Journal, 26*(4), 254-259. https://doi.org/10.1002/ffj.2054

Moriwaki, H., Tian, Y.-S., Kawashita, N., & Takagi, T. (2018). Mordred: A molecular descriptor calculator. *Journal of Cheminformatics, 10*, 4. https://doi.org/10.1186/s13321-018-0258-y

O'Boyle, N. M., Morley, C., & Hutchison, G. R. (2008). Pybel: a Python wrapper for the OpenBabel cheminformatics toolkit. *Chemistry Central Journal, 2*, 5. https://doi.org/10.1186/1752-153X-2-5

Okuyama, T., Miyashita, Y., Kanaya, S., Katsumi, H., Sasaki, S.-i., & Randić, M. (1988). Computer assisted structure-taste studies on sulfamates by pattern recognition method using graph theoretical invariants. *Journal of Computational Chemistry, 9*(6), 636-646. https://doi.org/10.1002/jcc.540090609

Oliveri, P. (2017). Class-modelling in food analytical chemistry: Development, sampling, optimisation and validation issues - A tutorial. *Analytica Chimica Acta, 982*, 9-19. https://doi.org/10.1016/j.aca.2017.05.013

Pallante, L., Korfiati, A., Androutsos, L., Stojceski, F., Bompotas, A., Giannikos, I., Raftopoulos, C., Malavolta, M., Grasso, G., Mavroudi, S., Kalogeras, A., Martos, V., Amoroso, D., Piga, D., Theofilatos, K., & Deriu, M. A. (2022). Toward a general and interpretable umami taste predictor using a multi-objective machine learning approach. *Scientific Reports, 12*, 21735. https://doi.org/10.1038/s41598-022-25935-3

Pieroni, A., Quave, C., Nebel, S., & Heinrich, M. (2002). Ethnopharmacy of the ethnic Albanians (Arbëreshë) of northern Basilicata, Italy. *Fitoterapia, 73*(3), 217-241. https://doi.org/10.1016/S0367-326X(02)00063-1

Pieroni, A., Houlihan, L., Ansari, N., Hussain, B., & Aslam, S. (2007). Medicinal perceptions of vegetables traditionally consumed by South-Asian migrants living in Bradford, Northern England. *Journal of Ethnopharmacology, 113*(1), 100-110. https://doi.org/10.1016/j.jep.2007.05.009

Rhyu, M.-R., Kim, Y., & Lyall, V. (2021). Interactions between Chemesthesis and Taste: Role of TRPA1 and TRPV1. *International Journal of Molecular Sciences, 22*(7), 3360. https://doi.org/10.3390/ijms22073360

Richard, A. M., & Williams, C. R. (2002). Distributed structure-searchable toxicity (DSSTox) public database network: A proposal. *Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis, 499*(1), 27-52. https://doi.org/10.1016/S0027-5107(01)00289-5

1313 Rodgers, S., Glen, R. C., & Bender, A. (2006). Characterizing bitterness: Identification of key structural features and development of a classification model. *Journal of Chemical Information and Modeling, 46*(2), 569-576. https://doi.org/10.1021/ci0504418

1316 Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling, 50*(5), 742-754. https://doi.org/10.1021/ci100050t

1318 Rojas, C., Duchowicz, P. R., Pis Diez, R., & Tripaldi, P. (2016a). Applications of quantitative structure-relative sweetness relationships in food chemistry. In A. G. Mercader, P. R. Duchowicz & P. M. Sivakumar (Eds.), *Chemometrics applications and research: QSAR in medicinal chemistry*, (pp. 317-339). Boca Raton (USA): Apple Academic Press.

1322 Rojas, C., Tripaldi, P., & Duchowicz, P. R. (2016b). A new QSPR study on relative sweetness. *International Journal of Quantitative Structure-Property Relationships, 1*(1), 78-92. https://doi.org/10.4018/IJQSPR.2016010104

1325 Rojas, C., Ballabio, D., Consonni, V., Tripaldi, P., Mauri, A., & Todeschini, R. (2016c). Quantitative structure-activity relationships to predict sweet and non-sweet tastes. *Theoretical Chemistry Accounts, 135*, 66. https://doi.org/10.1007/s00214-016-1812-1

1328 Rojas, C., Todeschini, R., Ballabio, D., Mauri, A., Consonni, V., Tripaldi, P., & Grisoni, F. (2017). A QSTR-based expert system to predict sweetness of molecules. *Frontiers in Chemistry, 5*, 53. https://doi.org/10.3389/fchem.2017.00053

1331 Rojas, C., Ballabio, D., Pacheco Sarmiento, K., Pacheco Jaramillo, E., Mendoza, M., & García, F. (2022). *ChemTastesDB*: A curated database of molecular tastants. *Food Chemistry: Molecular Sciences, 4*, 100090. https://doi.org/10.1016/j.fochms.2022.100090

1334 Roper, S. D. (2007). Signal transduction and information processing in mammalian taste buds. *Pflügers Archiv-European Journal of Physiology, 454*, 759-776. https://doi.org/10.1007/s00424-007-0247-x

1337 Ruddigkeit, L., & Reymond, J.-L. (2014). The chemical space of flavours. In K. Martinez-Mayorga & J. L. Medina-Franco (Eds.), *Foodinformatics: Applications of chemical information to food chemistry*, (pp. 83-96). Cham (Switzerland): Springer.

1340 Sahigara, F., Mansouri, K., Ballabio, D., Mauri, A., Consonni, V., & Todeschini, R. (2012). Comparison of different approaches to define the applicability domain of QSAR models. *Molecules, 17*(5), 4791-4810. https://doi.org/10.3390/molecules17054791

1343 Sayers, E. W., Beck, J., Bolton, E. E., Bourexis, D., Brister, J. R., Canese, K., Comeau, D. C., Funk, K., Kim, S., & Klimke, W. (2021). Database resources of the national center for biotechnology information. *Nucleic Acids Research, 49*(D1), D10-D17. https://doi.org/10.1093/nar/gkaa892

Schieberle, P., & Hofmann, T. (2016). Mapping the combinatorial code of food flavors by means of molecular sensory science approach. In H. Jeleń (Ed.), *Food flavors: Chemical, sensory and technological properties*, (pp. 413-438 ). Boca Raton (USA): CRC Press.

Schrödinger LLC. (2015). QikProp, New York, NY.

Schrödinger LLC. (2017). Canvas, New York, NY.

Shiyan, R., Liping, S., Xiaodong, S., Jinlun, H., & Yongliang, Z. (2021). Novel umami peptides from tilapia lower jaw and molecular docking to the taste receptor T1R1/T1R3. *Food Chemistry, 362*, 130249. https://doi.org/10.1016/j.foodchem.2021.130249

Spillane, W. J., & McGlinchey, G. (1981). Structure-activity studies on sulfamate sweeteners II: Semiquantitative structure-taste relationship for sulfamate ($RNHSO_3^-$) sweeteners-The role of R. *Journal of Pharmaceutical Sciences, 70*(8), 933-935. https://doi.org/10.1002/jps.2600700826

Spillane, W. J., McGlinchey, G., Muircheartaigh, I. Ó., & Benson, G. A. (1983). Structure-activity studies on sulfamate sweetners III: Structure-taste relationships for heterosulfamates. *Journal of Pharmaceutical Sciences, 72*(8), 852-856. https://doi.org/10.1002/jps.2600720804

Spillane, W. J., & Sheahan, M. B. (1989). Semi-quantitative and quantitative structure-taste relationships for carbo and hetero-sulphamate ($RNHSO_3^-$) sweeteners. *Journal of the Chemical Society, Perkin Transactions 2*(7), 741-746. https://doi.org/10.1039/P29890000741

Spillane, W. J., Sheahan, M. B., & Ryder, C. A. (1993). Synthesis and taste properties of sodium disubstituted phenylsulfamates. Structure-taste relationships for sweet and bitter/sweet sulfamates. *Food Chemistry, 47*(4), 363-369. https://doi.org/10.1016/0308-8146(93)90178-I

Spillane, W. J., Ryder, C. A., Walsh, M. R., Curran, P. J., Concagh, D. G., & Wall, S. N. (1996). Sulfamate sweeteners. *Food Chemistry, 56*(3), 255-261. https://doi.org/10.1016/0308-8146(96)00022-2

Spillane, W. J., Ryder, C. A., Curran, P. J., Wall, S. N., Kelly, L. M., Feeney, B. G., & Newell, J. (2000). Development of structure-taste relationships for sweet and non-sweet heterosulfamates. *Journal of the Chemical Society, Perkin Transactions 2, 2*, 1369-1374. https://doi.org/10.1039/B002482L

Spillane, W. J., Feeney, B. G., & Coyle, C. M. (2002). Further studies on the synthesis and tastes of monosubstituted benzenesulfamates. A semi-quantitative structure-taste relationship for the meta-compounds. *Food Chemistry, 79*(1), 15-22. https://doi.org/10.1016/S0308-8146(02)00169-3

Spillane, W. J., Kelly, L. M., Feeney, B. G., Drew, M. G., & Hattotuwagama, C. K. (2003). Synthesis of heterosulfamates. Search for structure-taste relationships. *Arkivoc, VII*, 297-309.
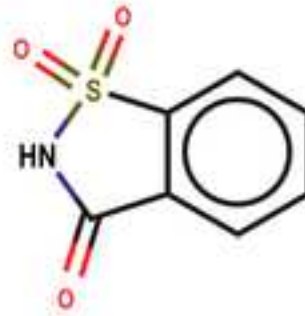
Spillane, W. J., Kelly, D. P., Curran, P. J., & Feeney, B. G. (2006). Structure-taste relationships for disubstituted phenylsulfamate tastants using classification and regression tree (CART) analysis. *Journal of Agricultural and Food Chemistry, 54*(16), 5996-6004. https://doi.org/10.1021/jf0606656

Spillane, W. J., Coyle, C. M., Feeney, B. G., & Thompson, E. F. (2009). Development of structure-taste relationships for thiazolyl-, benzothiazolyl-, and thiadiazolylsulfamates. *Journal of Agricultural and Food Chemistry, 57*(12), 5486-5493. https://doi.org/10.1021/jf9002472

Suárez-Estrella, D., Borgonovo, G., Buratti, S., Ferranti, P., Accardo, F., Pagani, M. A., & Marti, A. (2021). Sprouting of quinoa (Chenopodium quinoa Willd.): Effect on saponin content and relation to the taste and astringency assessed by electronic tongue. *LWT-Food Science and Technology, 144*, 111234. https://doi.org/10.1016/j.lwt.2021.111234

Suess, B., Festring, D., & Hofmann, T. (2015). Umami compounds and taste enhancers. In J. K. Parker, J. S. Elmore & L. Methven (Eds.), *Flavour development, analysis and perception in food and beverages*, (pp. 331-351). Cambridge (UK): Woodhead Publishing.

Takahashi, Y., Miyashita, Y., Tanaka, Y., Abe, H., & Sasaki, S. (1982). A consideration for structure-taste correlations of perillartines using pattern-recognition techniques. *Journal of Medicinal Chemistry, 25*(10), 1245-1248. https://doi.org/10.1021/jm00352a030

Takahashi, Y., Abe, H., Miyashita, Y., Tanaka, Y., Hayasaka, H., & Sasaki, S. I. (1984). Discriminative structural analysis using pattern recognition techniques in the structure-taste problem of perillartines. *Journal of Pharmaceutical Sciences, 73*(6), 737-741. https://doi.org/10.1002/jps.2600730608

Todeschini, R., & Consonni, V. (2009). *Molecular descriptors for chemoinformatics (2 volumes)* (Second ed.). Weinheim (Germany): Wiley-VCH.

Todeschini, R., Ballabio, D., Cassotti, M., & Consonni, V. (2015a). N3 and BNN: Two new similarity based classification methods in comparison with other classifiers. *Journal of Chemical Information and Modeling, 55*(11), 2365-2374. https://doi.org/10.1021/acs.jcim.5b00326

Todeschini, R., Ballabio, D., & Consonni, V. (2015b). Distances and other dissimilarity measures in chemometrics. In R. A. Meyers (Ed.), *Encyclopedia of analytical chemistry: Applications, theory and instrumentation*, (pp. 1-34): JohnWiley & Sons, Ltd.

Tuwani, R., Wadhwa, S., & Bagler, G. (2019). BitterSweet: Building machine learning models for predicting the bitter and sweet taste of small molecules. *Scientific Reports, 9*, 7155. https://doi.org/10.1038/s41598-019-43664-y

van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research, 9*, 2579-2605.

1414 Vapnik, V. (1998). The support vector method of function estimation. In J. A. K. Suykens & J.
1415      Vandewalle (Eds.), *Nonlinear modeling: Advanced black-box techniques*, (pp. 55-85). Boston
1416      (USA): Kluwer Academic Publishers.

1417 Walters, D. E. (2006). Analysing and predicting properties of sweet-tasting compounds. In W. J.
1418      Spillane (Ed.), *Optimising sweet taste in foods*, (pp. 283-291). Boca Raton (USA): CRC Press.

1419 Wiener, A., Shudler, M., Levit, A., & Niv, M. Y. (2012). BitterDB: A database of bitter compounds.
1420      *Nucleic Acids Research, 40*(D1), D413-D419. https://doi.org/10.1093/nar/gkr755

1421 Wold, S., & Eriksson, L. (1995). Statistical validation of QSAR results. Validation tools. In H. van
1422      de Waterbeemd (Ed.), *Chemometric methods in molecular design*, (pp. 309-318). Weinheim
1423      (Germany): VCH Publishers.

1424 Wong, D. W. (2018). *Mechanism and theory in food chemistry* (2nd ed.). Cham (Switzerland):
1425      Springer.

1426 Xiu, H., Liu, Y., Yang, H., Ren, H., Luo, B., Wang, Z., Shao, H., Wang, F., Zhang, J., & Wang, Y.
1427      (2022). Identification of novel umami molecules via QSAR models and molecular docking.
1428      *Food & Function, 13*, 7529-7539. https://doi.org/10.1039/D2FO00544A

1429 Yang, Z.-F., Xiao, R., Xiong, G.-L., Lin, Q.-L., Liang, Y., Zeng, W.-B., Dong, J., & Cao, D.-s. (2022).
1430      A novel multi-layer prediction approach for sweetness evaluation based on systematic
1431      machine learning modeling. *Food Chemistry, 372*, 131249.
1432      https://doi.org/10.1016/j.foodchem.2021.131249

1433 Yu, Z., Jiang, H., Guo, R., Yang, B., You, G., Zhao, M., & Liu, X. (2018). Taste, umami-enhance
1434      effect and amino acid sequence of peptides separated from silkworm pupa hydrolysate. *Food
1435      Research International, 108*, 144-150. https://doi.org/10.1016/j.foodres.2018.02.047

1436 Zheng, S., Jiang, M., Zhao, C., Zhu, R., Hu, Z., Xu, Y., & Lin, F. (2018). e-Bitter: Bitterant prediction
1437      by the consensus voting from the machine-learning methods. *Frontiers in Chemistry, 6*, 82.
1438      https://doi.org/10.3389/fchem.2018.00082

1439 Zheng, S., Chang, W., Xu, W., Xu, Y., & Lin, F. (2019). e-Sweet: A machine-learning based platform
1440      for the prediction of sweetener and its relative sweetness. *Frontiers in Chemistry, 7*, 35.
1441      https://doi.org/10.3389/fchem.2019.00035

1442

**Figure 1.** Taste changes of saccharin when introducing the nitro and amino molecular fragments in diverse position of the chemical scaffold.

**Figure 2.** Representation of classification boundaries (black lines) between sweet (blue) and bitter (red) chemicals in the space of the first two t-SNE dimensions (latent variables for PLSDA). The results are presented for different classifiers.

**Figure 3.** Number of molecules (expressed as log10) used for the calculation of models for taste prediction vs publication year.

Saccharin (sweet)

5-Nitrosaccharin (Bitter)

p-Nitrosaccharin (Sweet/Bitter)

6-Aminosaccharin (Sweet/Tasteless)

5-Aminosaccharin (Tasteless)

Figure 2

Figure 3

**Table 1.** Classification-based machine learning models for the discrimination between sweet and bitter tastants. $d$ is the number of descriptors, $n$ is the number of molecules.

**Table 2.** Classification-based machine learning models for the discrimination between sweet and non-sweet tastants. $d$ is the number of descriptors, $n$ is the number of molecules.

**Table 3.** Classification-based machine learning models for the prediction of bitterness. $d$ is the number of descriptors, $n$ is the number of molecules.

**Table 1**

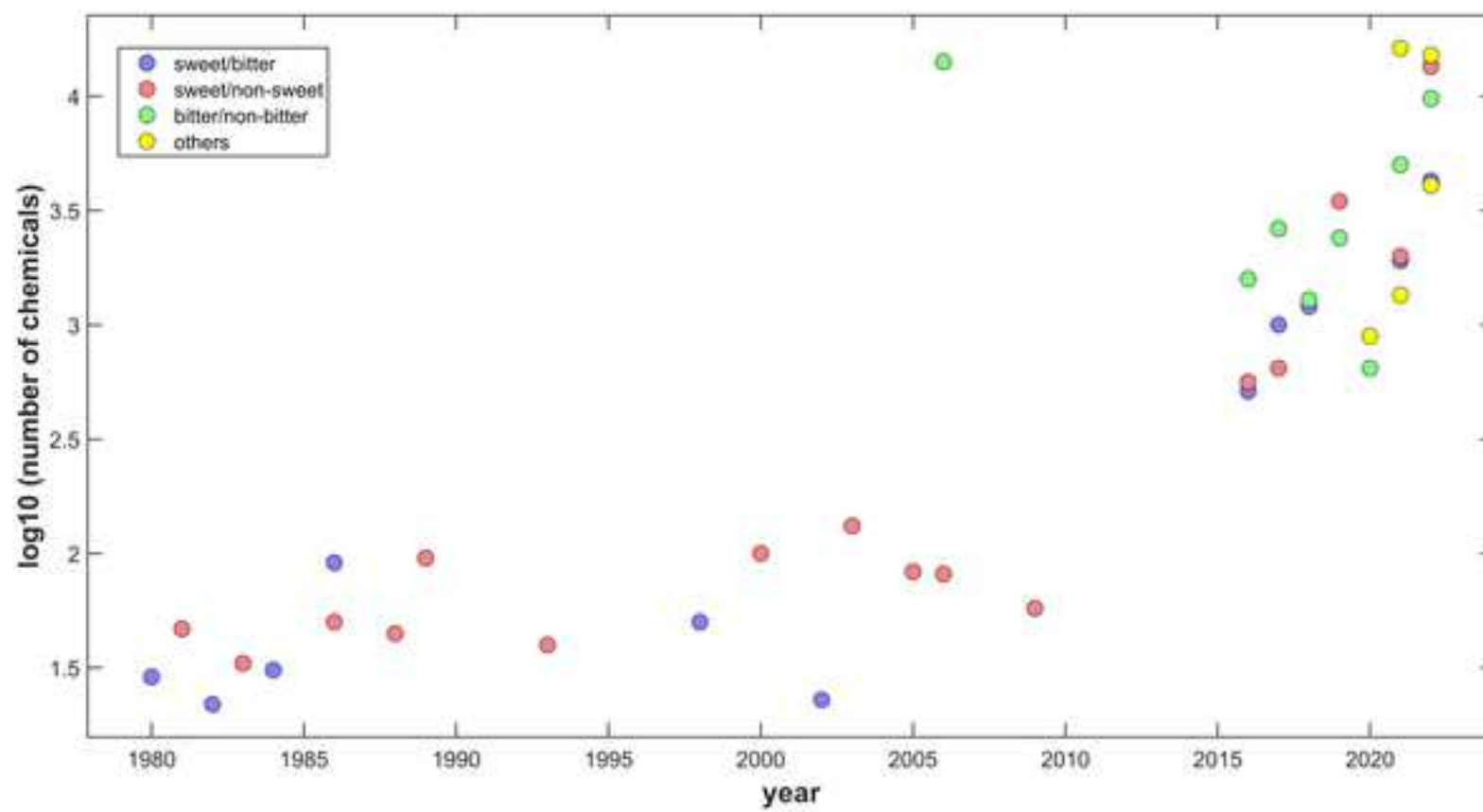| reference | URL | Model name | Classifier | $d$ | training | | | | test | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | $n$ | NER | AUC | F-score | $n$ | NER | AUC | F-score |
| (Kier, 1980) | -- | -- | LDA | 2 | 20 | 0.850 | -- | -- | 9 | 0.875 | -- | -- |
| (Takahashi et al., 1982) | -- | -- | kNN | 3 | 22 | 0.909 | -- | -- | -- | -- | -- | -- |
| | -- | | LDA | | | 1 | -- | -- | -- | -- | -- | -- |
| (Takahashi et al., 1984) | -- | -- | LDA | 3 | 22 | 1 | -- | -- | 9 | 0.775 | -- | -- |
| | | | | 2 | | 0.955 | -- | -- | | 0.775 | -- | -- |
| (Miyashita et al., 1986a) | -- | -- | SIMCA | 5 | 91 | 0.840 | | | -- | -- | -- | -- |
| (Drew et al., 1998) | -- | -- | DA | 11 | 50 | 1 | -- | -- | -- | -- | -- | -- |
| (Spillane et al., 2002) | -- | -- | Biplot | 2 | 23 | 0.862 | -- | -- | -- | -- | -- | -- |
| | | | LDA | 4 | | 0.850 | -- | -- | -- | -- | -- | -- |
| | | | QDA | 4 | | 0.900 | -- | -- | -- | -- | -- | -- |
| (Rojas et al., 2016c) | -- | -- | kNN | 4 | 356 | 0.864 | -- | -- | 152 | 0.789 | -- | -- |
| (Chéron et al., 2017) | http://sebfiorucci.free.fr/SweetenersDB/ | -- | RF | 5[a] | 796 | 0.997 | -- | -- | 200 | 0.914 | -- | -- |
| (Banerjee & Preissner, 2018) | -- | BitterSweetForest | RF | 2,048 | 961 | 0.950[b] | 0.980 | 0.940 | 241 | 0.967[b] | 0.980 | 0.920 |
| (Goel et al., 2021) | -- | -- | RF | 8 | 1,537 | 0.908 | -- | -- | 385 | 0.855 | -- | -- |
| (Bo et al., 2022) | -- | BitterSweetMLP-Fingerprint | MLP | 17 | 1,637 | 0.870 | 0.950 | -- | 409 | 0.880 | 0.950 | -- |
| (Maroni et al., 2022) | https://github.com/gabribg88/VirtuousSweetBitter https://virtuoush2020.com/ | -- | GBM | 9 | 2,195 | 0.893 | 0.950 | 0.883 | -- | -- | -- | -- |

[a] number of descriptors for the tree depth; [b] calculated as Accuracy (ACC)

**Table 2**

| reference | URL | Model name | Classifier | d | training | | | | test | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | n | NER | AUC | F-score | n | NER | AUC | F-score |
| (Spillane & McGlinchey, 1981) | -- | -- | DA-plot | 2 | 47 | 0.957[b] | -- | -- | -- | -- | -- | -- |
| (Spillane et al., 1983) | -- | -- | LDA | 3 | 33 | 0.807 | -- | -- | -- | -- | -- | -- |
| (Miyashita et al., 1986b) | -- | -- | SIMCA | 4 | 50 | 0.798 | -- | -- | -- | -- | -- | -- |
| (Okuyama et al., 1988) | -- | -- | SIMCA | 1[a] | 25 | 0.868 | -- | -- | -- | -- | -- | -- |
| | | | | | 20 | 0.808 | -- | -- | -- | -- | -- | -- |
| (Spillane & Sheahan, 1989) | -- | -- | DA-plot | 2 | 17 | 0.824[b] | -- | -- | -- | -- | -- | -- |
| | | | LDA | 3 | 23 | 0.642 | -- | -- | -- | -- | -- | -- |
| | | | | | 56 | 0.773 | -- | -- | -- | -- | -- | -- |
| (Spillane et al., 1993) | -- | -- | DA-plot | 2 | 40 | -- | -- | -- | -- | -- | -- | -- |
| (Spillane et al., 2000) | -- | -- | QDA | 4 | 101 | 0.801 | -- | -- | -- | -- | -- | -- |
| (Spillane et al., 2003) | -- | -- | CART | 4 | 132 | 0.815 | -- | -- | -- | -- | -- | -- |
| (Kelly et al., 2005) | -- | -- | CART | 6 | 75 | 0.768 | -- | -- | 8 | 0.750 | -- | -- |
| (Spillane et al., 2006) | -- | -- | CART | 7 | 70 | 0.807 | -- | -- | 12 | 0.909 | -- | -- |
| (Spillane et al., 2009) | -- | -- | CART | 6 | 48 | 0.950 | -- | -- | 10 | 0.625 | -- | -- |
| (Rojas et al., 2016c) | -- | -- | kNN | 9 | 396 | 0.838 | -- | -- | 170 | 0.752 | -- | -- |
| (Rojas et al., 2017) | -- | -- | Expert System | -- | 488 | 0.892 | -- | -- | 161 | 0.848 | -- | -- |
| (Zheng et al., 2019) | https://www.dropbox.com/sh/1fmlv7nf6wofgcp/AADBJzFbbbiNRJUP0806wSyna?dl=0 | e-Sweet | Consensus | -- | 883 | 0.870 | -- | 0.850 | 221 | 0.900 | -- | 0.878 |
| (Tuwani et al., 2019) | https://github.com/cosylabiiit/bittersweet/ <br> https://cosylab.iiitd.edu.in/bittersweet/ | BitterSweet | AdaBoost | -- | 2,205 | 0.856 | 0.918 | 0.858 | 161 | 0.834 | 0.883 | 0.856 |
| (Fritz et al., 2021) | http://virtualtaste.charite.de/VirtualTaste/ | VirtualSweet | RF | -- | 1,608 | 0.970 | 0.990 | 0.870 | 403 | 0.893 | 0.951 | 0.888 |
| (Yang et al., 2022) | | -- | RF | 241 | 959 | 0.873 | 0.958 | -- | 241 | 0.920 | 0.971 | -- |

| reference | URL | Model name | Classifier | d | n | NER | AUC | F-score | n | NER | AUC | F-score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | https://github.com/ifyoungnet/ChemSweet | | XGBoost | 95 | 366 | 0.905 | 0.956 | -- | 92 | 0.926 | 0.974 | -- |
| | | | | 105 | 1,327 | 0.834 | 0.926 | -- | 333 | 0.841 | 0.920 | -- |
| | | | | 124 | 2,104 | 0.870 | 0.947 | -- | 527 | 0.867 | 0.947 | -- |
| | | | | 102 | 394 | 0.893 | 0.937 | -- | 100 | 0.876 | 0.956 | -- |
| | | | | 122 | 2,091 | 0.875 | 0.949 | -- | 522 | 0.889 | 0.961 | |
| (Bo et al., 2022) | -- | SweetMLP-Fingerprint | MLP | -- | 1,776 | 0.860 | 0.930 | -- | 444 | 0.900 | 0.940 | -- |
| | | SweetCNN | CNN | | | 0.860 | 0.900 | -- | | 0.850 | 0.900 | -- |
| (Lee et al., 2022) | -- | BoostSweet | Soft-vote consensus | -- | 1,832 | -- | -- | -- | 459 | 0.899 | 0.961 | 0.907 |

[a] number of principal components (PCs); [b] calculated as Accuracy (*ACC*)

**Table 3**

| reference | URL | Model name | Classifier | d | training | | | | test | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | n | NER | AUC | F-score | n | NER | AUC | F-score |
| (Rodgers et al., 2006) | -- | -- | Naïve Bayes | 10 | 14,179 | 0.805 | -- | -- | -- | -- | -- | -- |
| (Huang et al., 2016) | http://mdl.shsmu.edu.cn/BitterX | BitterX | SVM | 46 | 862 | 0.879[b] | -- | -- | 216 | 0.915[b] | 0.950 | -- |
| | | | | 35 | 416 | 0.767[b] | -- | -- | 104 | 0.798[b] | 0.823 | -- |
| (Dagan-Wiener et al., 2017) | https://github.com/Niv-Lab/BitterPredict1 | BitterPredict | AdaBoost | 16[a] | 1,827 | 0.921 | -- | -- | 781 | 0.812 | -- | -- |
| (Zheng et al., 2018) | https://www.dropbox.com/sh/3sebvza3qzmazda/AADgpCRXJtHAJzS8DK_P-q0ka?dl=0 | e-Bitter | Consensus | -- | 1,040 | -- | -- | -- | 259 | 0.929[b] | -- | 0.936 |
| (Tuwani et al., 2019) | https://github.com/cosylabiiit/bittersweet/ https://cosylab.iiitd.edu.in/bittersweet/ | BitterSweet | RF | -- | 2,257 | 0.754 | 0.852 | 0.698 | 154 | 0.819 | 0.880 | 0.838 |

| Reference | URL | Tool | Method | Desc.[a] | N | ACC/M1 | M2 | M3 | N | ACC/M1 | M2 | M3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| (Charoenkwan et al., 2020a) | http://camt.pythonanywhere.com/ | iBitter-SCM | SCM | -- | 512 | 0.871[b] | | | 128 | 0.844[b] | | |
| (Margulis et al., 2021) | -- | BitterIntense | XGBoost | 8 | 616 | 0.870[b] | -- | 0.820 | 105 | 0.790 | -- | 0.700 |
| (Charoenkwan et al., 2021a) | http://pmlab.pythonanywhere.com/BERT4Bitter | BERT4Bitter | BERT | -- | 512 | 0.861[b] | 0.915 | -- | 128 | 0.922[b] | 0.964 | -- |
| (Fritz et al., 2021) | http://virtualtaste.charite.de/VirtualTaste/ | VirtualBitter | RF | -- | 1,289 | 0.960 | 0.975 | 0.946 | 323 | 0.898 | 0.956 | 0.882 |
| (Charoenkwan et al., 2021b) | http://camt.pythonanywhere.com/iBitter-Fuse | iBitter-Fuse | SVM | 36 | 512 | 0.918[b] | 0.937 | -- | 128 | 0.930[b] | 0.933 | -- |
| (Bai et al., 2021) | -- | CBDPS | XGBoost | -- | 1,296 | 0.882[b] | -- | 0.881 | 112 | -- | -- | -- |
| (Bo et al., 2022) | -- | BitterMLP-Descriptor | MLP | 15 | 1,787 | 0.830 | 0.920 | -- | 446 | 0.820 | 0.940 | -- |
| | | BitterCNN | CNN | -- | | 0.770 | 0.870 | -- | | 0.790 | 0.880 | -- |
| (Margulis et al., 2022) | https://github.com/YuliSl/BitterMatch | BitterMatch | XGBoost | 20 | 3,601 | 0.759[c] | -- | -- | 900 | -- | -- | -- |
| | | | | | 242 | 0.699[c] | -- | -- | 61 | -- | -- | -- |
| (De León et al., 2022) | -- | Premexotac | SVM | 512 | 2,272 | 0.836[b] | -- | -- | 568 | 0.870[b] | -- | -- |
| | | | AdaBoost | 18 | | 0.842[b] | -- | -- | | 0.847[b] | -- | -- |

[a] descriptors with the most significant contribution; [b] calculated as Accuracy (*ACC*), [c] reported as recall-precision