



UNIVERSITAT AUTÒNOMA DE BARCELONA

---

# Notes: AI application for azophotoswitches' optimization with pharmacological interest

---

AUTHOR : SERGIO CASTAÑEIRAS MORALES  
SUPERVISOR : MIQUEL MORENO FERRER  
CO-SUPERVISOR : ÀNGELS GONZALEZ LAFONT

FINAL DEGREE PROJECT  
BACHELOR'S DEGREE IN CHEMISTRY

2024-2025

*“The dumbest people I know are those who know it all.”*

**Malcolm S. Forbes**

# 1 Definitions

**Definition 1.**  $IC_{50}$ : Half maximal inhibitory concentration "IC<sub>50</sub> is the concentration of drug required for 50% inhibition. IC<sub>50</sub> is an operational term dependent on the assay conditions. IC<sub>90</sub> or IC<sub>99</sub> is sometimes used when complete inhibition is required. Calculation of the fractional occupancy shows that IC<sub>90</sub> concentration is approximately 10-fold greater than the IC<sub>50</sub> concentration assuming one-site binding at equilibrium with a Hill coefficient of 1."<sup>1</sup> For this project, we aim to identify substances with the lowest possible IC<sub>50</sub>, as our goal is to minimize the presence of foreign substances in the living organism.

**Definition 2.** Molecular descriptor: "A molecular descriptor is the final result of a logical and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into a useful number or the result of some standardized experiment."<sup>2</sup>

**Definition 3.** Constitutional Indices Descriptors: Kind of descriptor based on the constitutional composition of a molecule. They describe the basic composition of the molecule without considering its geometry or connectivity of it.

**Definition 4.** Ring Descriptors: Kind of descriptor that describes the ring systems of our molecule. Quantifies properties such as the ring connectivity, type of ring (odd/even ring), aromaticity & saturation.

**Definition 5.** Topological Indices: These descriptors are based on the connectivity of atoms in a molecule without considering three-dimensional coordinates in the context of chemical graph theory. They help to describe the molecular graph (atoms as vertices, bonds as edges) using metrics like the Wiener index or the Zagreb index.

**Definition 6.** Wiener index: Specific topological descriptor based defined as the sum of the

lengths of the shortest paths between all possible vertices in the chemical graph. Mathematically speaking, if we have a molecule with a set of  $\{A_i\}_{i=0}^n$  atoms and we denote  $d_B(A_i, A_j)$  the bond distance between atoms  $A_i$  and  $A_j$ , we define the Wiener index of this molecule as,

$$Ind_{Wiener} = \sum_{i=0}^n \sum_{j>i}^n d_B(A_i, A_j) \quad (1)$$

For instance, if we take the n-Butane & Isobutane chemical graph from the Figure (1) the Wiener index would be,

$$d_B(C_1, C_2) + d_B(C_1, C_3) + d_B(C_1, C_4) \quad (2)$$

$$+ d_B(C_2, C_3) + d_B(C_2, C_4) + d_B(C_3, C_4) \quad (3)$$

Though this quantity is computed equally for both, n-Butane & Isobutane, it results different due to their connectivity. For the n-Butane:

$$1 + 2 + 3 + 1 + 2 + 1 = 10 \quad (4)$$

While for the Isobutane:

$$1 + 1 + 1 + 2 + 2 + 2 = 9 \quad (5)$$

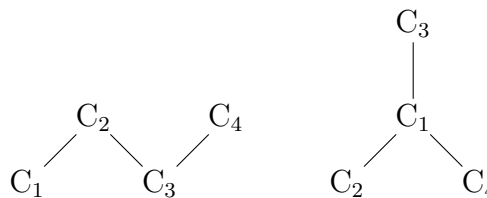


Figure 1: n-Butane's & Isobutane's chemical graphs.

**Definition 7.** Zagreb index: "For a (molecular) graph, the first Zagreb index  $M_1$  is equal to the sum of squares of the degrees of vertices, and the second Zagreb index  $M_2$  is equal to the sum of the products of the degrees of pairs of adjacent vertices."<sup>3</sup>

**Definition 8.** *Walk and Path Counts Descriptors:* Kind of descriptor focused on counting specific paths or walks (connected sequences of atoms) in the molecular graph. For instance we can get a Walk and Path Counts Descriptor quantifying the number of Eulerian paths in a molecule's graph. They give insights into the molecule's connectivity and complexity.

**Definition 9.** *Connectivity Indices:* Connectivity Indices are a class of molecular descriptors that quantify the connectivity or bonding patterns between atoms in a molecule based on its topology (i.e., the molecular graph). These indices capture the degree to which atoms are connected, considering how atoms are bonded to each other rather than their physical 3D arrangement. They are typically derived from graph theory, where atoms are represented as vertices and bonds as edges.<sup>4</sup>

**Definition 10.** *Randic Index:* Specific connectivity index defined as the following,

$$\sum_{\text{all bonds}} \frac{1}{\sqrt{d_i d_j}} \quad (6)$$

where  $d_i$  and  $d_j$  denote the number of bonds of atoms  $i$  and  $j$  connected by a bond.<sup>5</sup>

**Definition 11.** *Information Indices:* Kind of descriptor based on the information theory, where the molecular structure is represented as a distribution of information (or entropy). They measure molecular complexity and diversity.

**Definition 12.** *2D-Matrix Based Descriptors:* Kind of descriptor based in the matricidal representation of a molecule.<sup>6</sup>

**Definition 13.** *Adjacency Matrix:* Specific 2D-Matrix Based Descriptor based on the matricidal representation of a molecule. If we have  $\{x\}_{i=1}^n$  atoms that form our molecule, labeled from 1 to  $n$ , the adjacent matrix  $A$  will contain as the element  $a_{i,j}$  the number 1 if the atom  $i$

and  $j$  are connected and 0 otherwise. (We can also naturally add the possibilities of  $a_{i,j} = 2, 3$  for double and triple bonds).

For instance in the Figure (1) the adjacent matrix for the  $n$ -Butane is

$$n\text{-Butane } AM = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix} \quad (7)$$

Trivially we observe that the main diagonal will always be filled with 0s and the matrices will also be doubly symmetrical (respect both diagonals) due to the bounds transitivity and reflectivities properties.

**Definition 14.** *Distance Matrix:* Specific 2D-Matrix Based Descriptor based on the computation of the shortest pathways between atoms. If we have  $\{x\}_{i=1}^n$  atoms that form our molecule, labeled from 1 to  $n$ , the adjacent matrix  $A$  will contain as the element  $a_{i,j}$  the shortest bond distance between atom  $i$  and atom  $j$ .

For instance in the Figure (1) the distance matrix for the  $n$ -Butane is

$$n\text{-Butane } DM = \begin{pmatrix} 0 & 1 & 2 & 3 \\ 1 & 0 & 1 & 2 \\ 2 & 1 & 0 & 1 \\ 3 & 2 & 1 & 0 \end{pmatrix} \quad (8)$$

**Definition 15.** *2D Autocorrelations:* a type of molecular descriptor that quantifies the correlations between properties of atoms within a molecule based on their positions relative to one another in the molecular structure. These descriptors analyse how specific atomic properties are distributed across the molecule, capturing the spatial relationships of these properties.

For a given property and distance  $d$ , the autocorrelation is calculated by summing the products of the property values for all pairs of atoms that are separated by that distance. Mathematically, the autocorrelation  $A_d(P)$  at distance  $d$

for a property  $P$  is given by,

$$A_d(P) = \sum_{d(i,j)=d} P(i)P(j) \quad (9)$$

**Definition 16.** *Burden Eigenvalues:* Kind of descriptor extracted from the Burden Matrix of a molecule,<sup>7</sup> which is a modification of the adjacency matrix from the definition (13).

**Definition 17.** *P-VSA-like Descriptors:* (Property-weighted Van der Waals Surface Area descriptors) descriptors deviated from the computation of the Van der Waals Surface of a molecule. For instance, we could talk about polarity indices, high/low electron density regions...

**Definition 18.** *ETA Indices:* (Extended Topochemical Atom) kind of indices that capture aspects of molecular topology in conjunction with atomic properties. They are often used in the analysis of electronic properties or biological activities.<sup>2</sup>

**Definition 19.** *Edge Adjacency Indices:* Kind of descriptor that relies on the chemical graph theory and is related to the number of bonds that each atom generates. This descriptors are tightly related to the molecule's complexity.<sup>2</sup>

**Definition 20.** *Fractional Group Counts (FGC):* kind of molecular descriptor related to the chemical functional groups. For instance, in the case of cyclopentanone Figure (2), the Ketones Group Counts is  $\frac{1}{6}$ , since we have 1 ketone for 6 atoms (we do not typically compute the hydrogen atoms).<sup>2</sup>

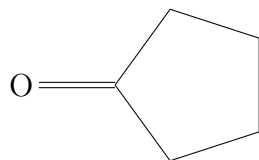


Figure 2: Cyclopentanone's chemical graph.

**Definition 21.** *Atom-Centered Fragments (ACF):* Kind of molecular descriptor concerned

about atoms separately. For a given molecule conformed by  $n$  atoms, and a specific Atom-Centered Fragment descriptor, we can compute  $n$  descriptors of this kind based one per each atom separately. Each descriptor describes an atom itself, its surroundings and its bounding types.<sup>2</sup>

**Definition 22.** *Atom-Type E-State Indices:* Kind of molecular descriptor that describes the electronic state of an individual atom. Some Atom-Type E-State Indices describe the hybridisation state (e.g.,  $sp^2$ ,  $sp^3$ ), the local bonding environment (e.g., single, double, or aromatic bonds), the presence of electronegative atoms nearby or the topological distance to other atoms in the molecule. For a given atom, a atom-type e-etate Index is the combination of 2 terms. The intrinsic term that depends on the atom hybridization and the perturbation term which is adjusted depending in the atom's surroundings.<sup>2</sup>

**Definition 23.** *Pharmacophore Descriptors:* Kind of descriptor referred to the functional groups of a given molecule and its molecular features. they describe properties such as hydrogen donation (-OH), hydrogen acceptation ( $C=O$ ), Positive or negative ionizable groups ( $-COO^-$ )... This kind of descriptor does not worry about the molecular structure rather than the molecular properties as a whole.<sup>8</sup>

**Definition 24.** *2D Atom Pairs:* Kind of descriptor that focuses on pairs of atoms and captures information about the types of atoms involved, as well as the distance between them in a two-dimensional representation of the molecular structure.

**Definition 25.** *Charge descriptors:* Kind of descriptor that describes a molecular charge, as a complete element as well as the composition of several elements. One of the key features of this kind of descriptors is their high dependency on dipole moment as well as the electronegativity of the atoms that it describes.

**Definition 26.** *Molecular Properties:* The most intuitive kind of descriptor. These descriptors are the traditional features of a substance such as its molecular weight, boiling point, solubility, density...

**Definition 27.** *Drug-Like Indices:* These descriptors quantify how closely the molecule follows the rules for drug-likeness, such as Lipinski's Rule of Five, which includes factors like molecular weight, hydrophobicity, and hydrogen bonding.

**Definition 28.** *MDE Descriptors (Molecular Distance Edge descriptors):* These describe the distances between atoms in a molecule based on the molecular graph. They are useful for capturing molecular shape and spatial relationships.

**Definition 29.** *Chirality Descriptors:* These quantify the chirality (handedness) of molecules. They reflect the 3D arrangement of atoms around chiral centers, which can influence how molecules interact with biological systems.

**Definition 30.** *QSAR: "Quantitative Structure-Activity Relationships are the final result of the process that starts with a suitable description of molecular structures and ends with some inference, hypothesis, and prediction on the behaviour of molecules in environmental, biological, and physico-chemical systems in analysis.*

QSARs are based on the assumption that the structure of a molecule (for example, its geometric, steric, and electronic properties) must contain features responsible for its physical, chemical, and biological properties and on the ability to capture these features into one or more numerical descriptors. By QSAR models, the biological activity (or property, reactivity, etc.) of a new designed or untested chemical can be inferred from the molecular structure of similar compounds whose activities (properties, reactivities, etc.) have already been assessed."<sup>6</sup>

**Definition 31.** *Machine Learning Model:* For a given problem, Machine Learning Model is the procedure to resolve the problem in a Machine Learning like format.

**Definition 32.** *The efficiency of a Machine Learning Model relies in the following quantities:*

- *True Positives (TP):* number of true positives results.
- *False Positives (FP):* number of false positives results.
- *True Negatives (TN):* number of true negatives results.
- *Fales Negatives (FN):* number of false negatives results.

**Definition 33.** *True Positive Rate:* quantity related to a Machine Learning Model's sensitivity defined as:

$$\frac{TP}{TP + FN} \quad (10)$$

**Definition 34.** *True Negative Rate:* quantity related to a Machine Learning Model's specificity defined as:

$$\frac{TN}{TN + FN} \quad (11)$$

**Definition 35.** *Classification Accuracy:* quantity related to a Machine Learning Model's effectiveness defined as:

$$\frac{TP + TN}{TP + FP + TN + FN} \quad (12)$$

**Definition 36.** *Matthews Correlation Coefficient:* quantity related to a Machine Learning Model's prediction capacity defined as:

$$\frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \quad (13)$$

A Matthews Correlation Coefficient equal to 1 stands for a perfect prediction a Matthews Correlation Coefficient equal to 0 indicates that the classifier's predictions are no better than random guessing, and a Matthews Correlation Coefficient equal to -1 stand for a total disagreement between predictions and actual outcomes.

**Definition 37.** Cyclooxygenase-2: (COX-2), also known as prostaglandin-endoperoxide synthase 2 is "An enzyme that speeds up the forma-

tion of substances that cause inflammation and pain. It may also cause tumor cells to grow. Some tumors have high levels of COX-2 and blocking its activity may reduce tumor growth. Also called cyclooxygenase-2 and prostaglandin-endoperoxide synthase 2."<sup>9</sup> This enzyme will be our main object of study. It is encoded by the PTGS2 gene and its main functionality is converting arachidonic acid to prostaglandin H<sub>2</sub> Figure (3).

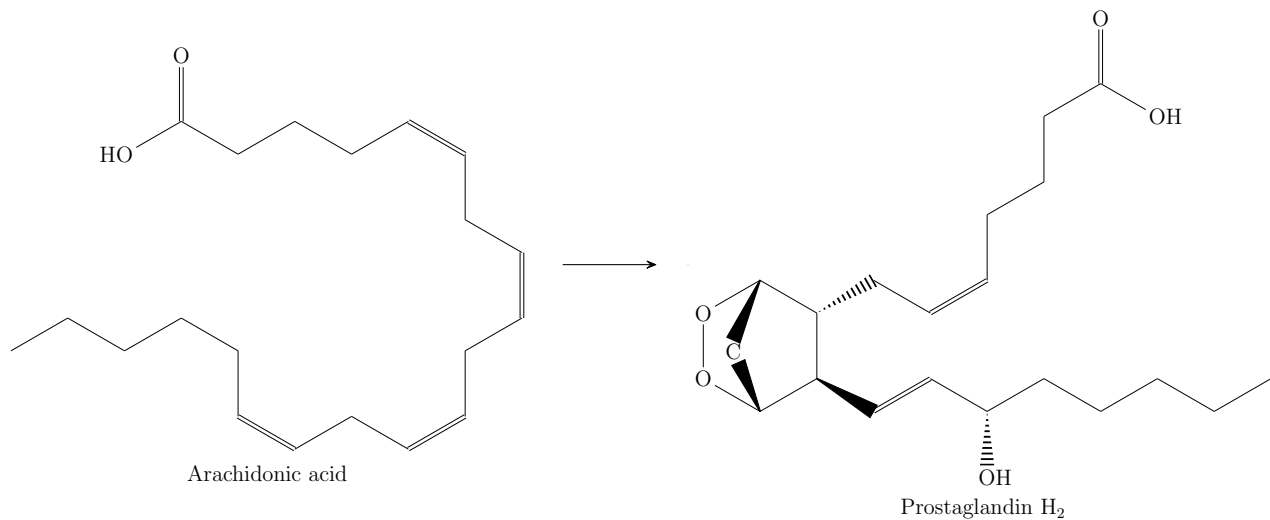


Figure 3: Reaction catalysed by the COX-2 between arachidonic acid to prostaglandin H<sub>2</sub>.

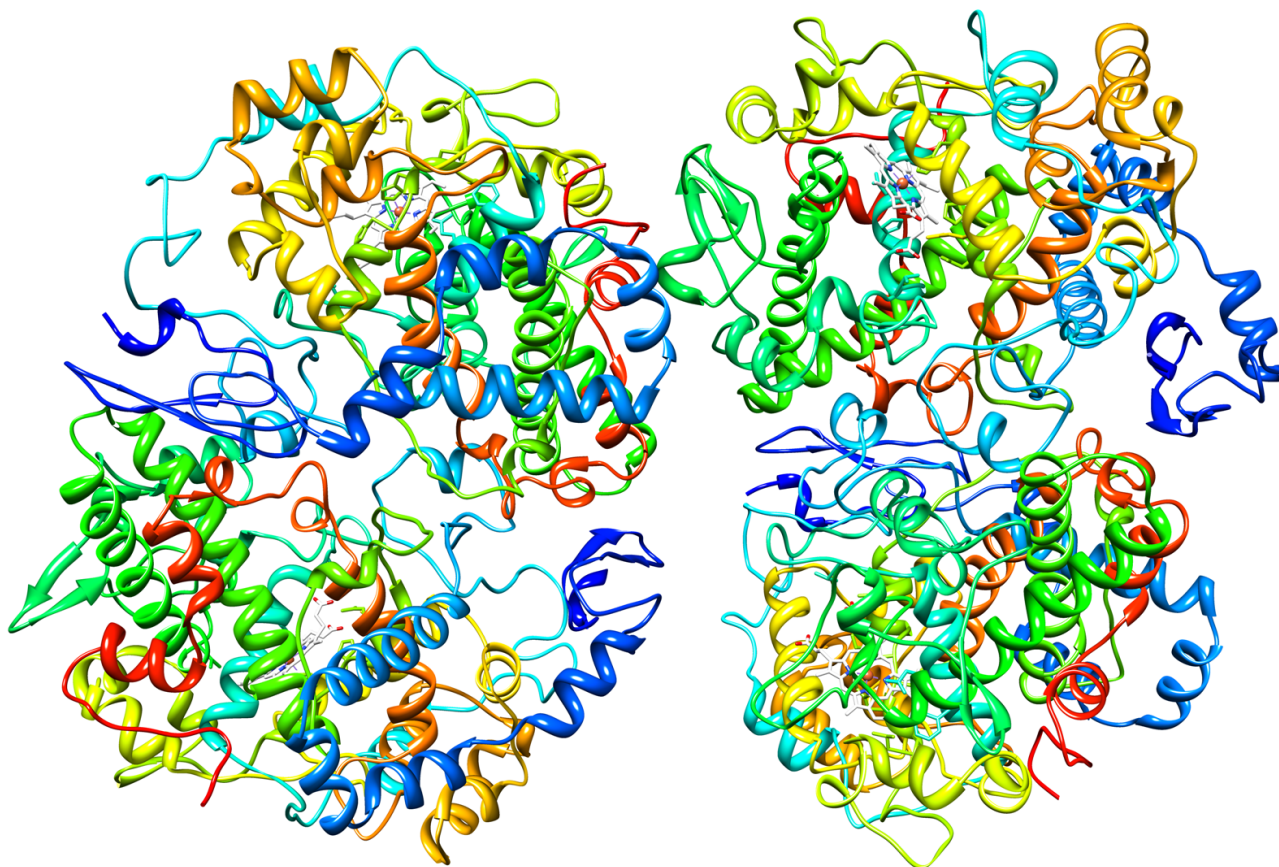


Figure 4: Crystal structure of uninhibited cyclooxygenase-2 Molecular graphics images were produced using the UCSF Chimera package.<sup>10</sup>

**Definition 38.** *Non-steroidal anti-inflammatory drug: (NSAID) kind of drug commonly used for therapeutic purposes that focuses in pain, inflammation and fever reduction and blood clots prevention. They are characterised by inhibiting the activity of enzymes called cyclooxygenase (COX). This drug type is known to produce side-effects such as:*

- NSAIDs block COX-1, which helps protect the stomach lining. Long-term use can lead to ulcers and gastrointestinal bleeding.
- Stomach pain, indigestion, heartburn.

- Particularly **COX-2 inhibitors**, may increase the risk of **heart attack**, stroke, and blood clots.
- Increased Blood Pressure.
- Allergic Reactions.
- Liver Toxicity.
- Fluid Retention and Edema.

Some well-known NSAIDs are Ibuprofen, Aspirin and Celecoxib where the last one is specific a specific inhibitor for the COX-2.<sup>11</sup>

In the Figure (5) we can observe the binding site of COX-2 for the celecoxib.



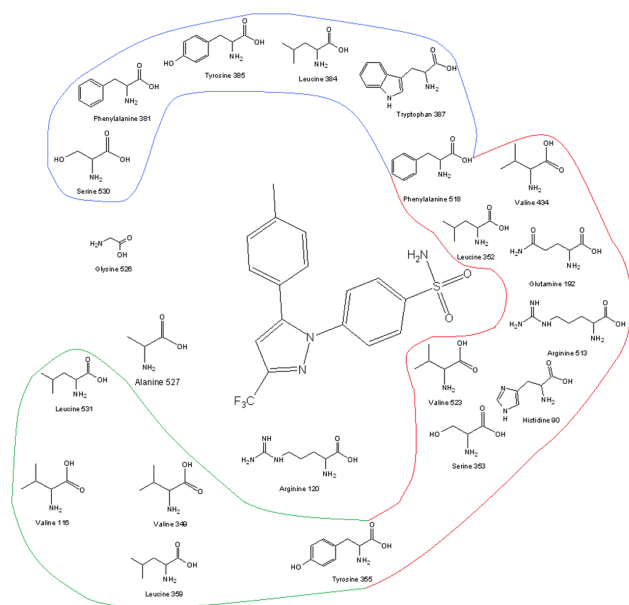


Figure 5: COX-2 receptor site and its amino acid profile with celecoxib in the binding site.<sup>12</sup>

**Definition 39.** *Celecoxib:*<sup>1</sup> drug known to be a selective COX-2 inhibitor (currently is not highly selective respect to newer drugs), see Figure (6). We will use the chemical descriptors of this molecule to compute its theoretical  $IC_{50}$  (one that our Machine Learning Model predicts) to have some sort of feed back of the ML reliability.

**Definition 40.** *Rofecoxib:*<sup>2</sup> drug known to be a selective COX-2 inhibitor, see Figure (7). We will use the chemical descriptors of this molecule to compute its theoretical  $IC_{50}$  (one that our Machine Learning Model predicts) to have some sort of feed back of the ML reliability.

**Definition 41.** *Ibuprofen:*<sup>3</sup> drug known to be a selective COX-1 and COX-2 inhibitor, see Figure (8). We will use the chemical descriptors of this molecule to compute its theoretical  $IC_{50}$  (one that our Machine Learning Model

predicts) to have some sort of feed back of the ML reliability.

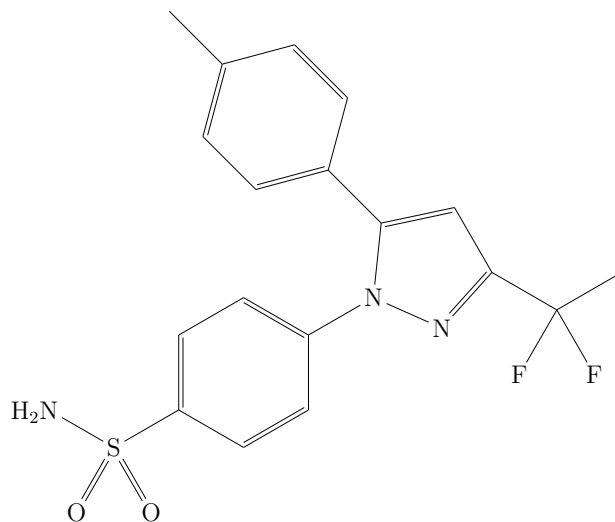


Figure 6: Chemical graph of Celecoxib.

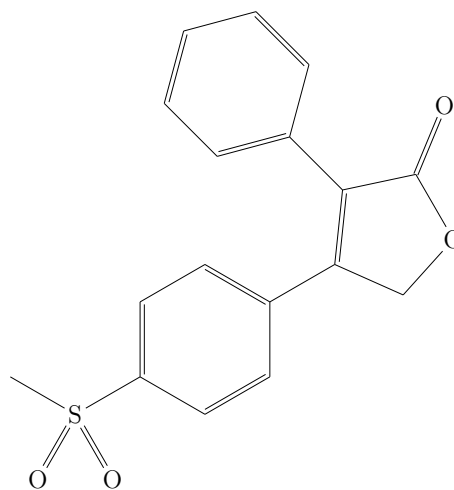


Figure 7: Chemical graph of Rofecoxib.

<sup>1</sup>UPAC name: 4-[5-(4-Methylphenyl)-3-(trifluoromethyl)pyrazol-1-yl]benzenesulfonamide

<sup>2</sup>UPAC name: 4-(4-methylsulfonylphenyl)-3-phenyl-5H-furan-2-one

<sup>3</sup>UPAC name: (RS)-2-(4-(2-Methylpropyl)phenyl)propanoic acid

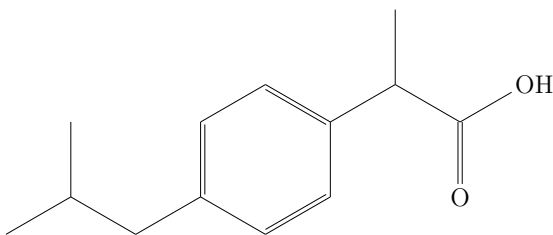


Figure 8: Chemical graph of Ibuprofen.

**Definition 42.** *Pearson correlation coefficient:* Given set of pairs of data  $\{(x_i, y_i)\}_{i=1}^n$  the pearson correlation factor  $r_{xy}$  is defined as,

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}, \quad (14)$$

where  $\bar{x}$  and  $\bar{y}$  stand for the average value of  $x_{i=1}^n$  and  $y_{i=1}^n$  respectively. Note that  $r_{xy} \in [-1, 1]$ . Therefore the sign of  $r_{xy}$  is tightly related to the sign of a linear regression, more precisely if  $x > 0$ , "y" generally<sup>4</sup> increases when "x" increases, as well as if  $x < 0$ , "y" decreases when "x" increases.

**Definition 43.** A decision tree is a classification algorithm based on a series of ordered if statements. The algorithm begins at the top of the tree, where a question is posed to the data. Depending on the answer, the data follows different branches, each corresponding to a subsequent question. This process is repeated at each node until the data reaches the bottom of the tree, where the path it has followed determines the classification of the given data.

For example, a simple depiction of a decision tree is shown in Figure (9), where the decision of whether to enjoy a golfing day or not is determined based on a series of questions.

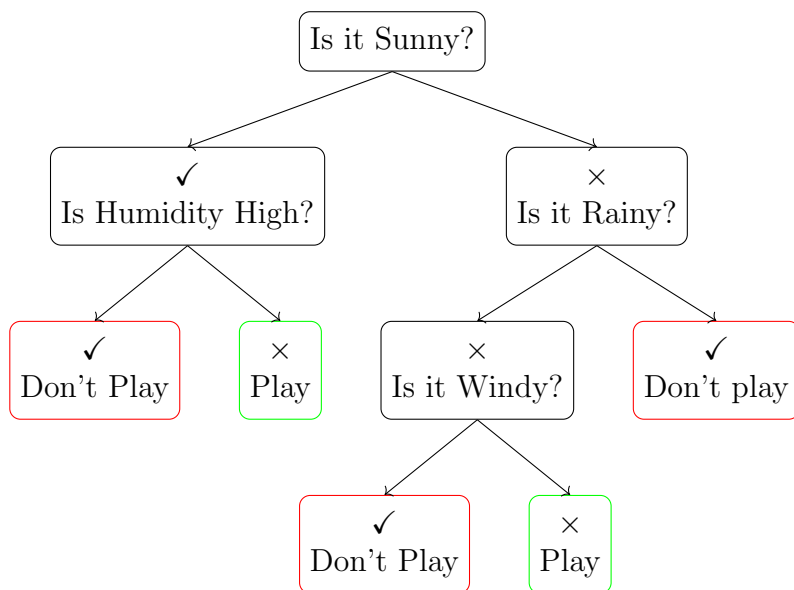


Figure 9: Example of a decision tree

This kind of algorithm is vividly present in the chemical landscape, for instance in the spectroscopy realm the determination of a molecule's symmetry group is provided by a decision tree such as Figure (10)

<sup>4</sup>We would like to remark that the word "generally" stands for "the majority of the cases", since "generally" is commonly interpreted as a non-scientific/non-objective word

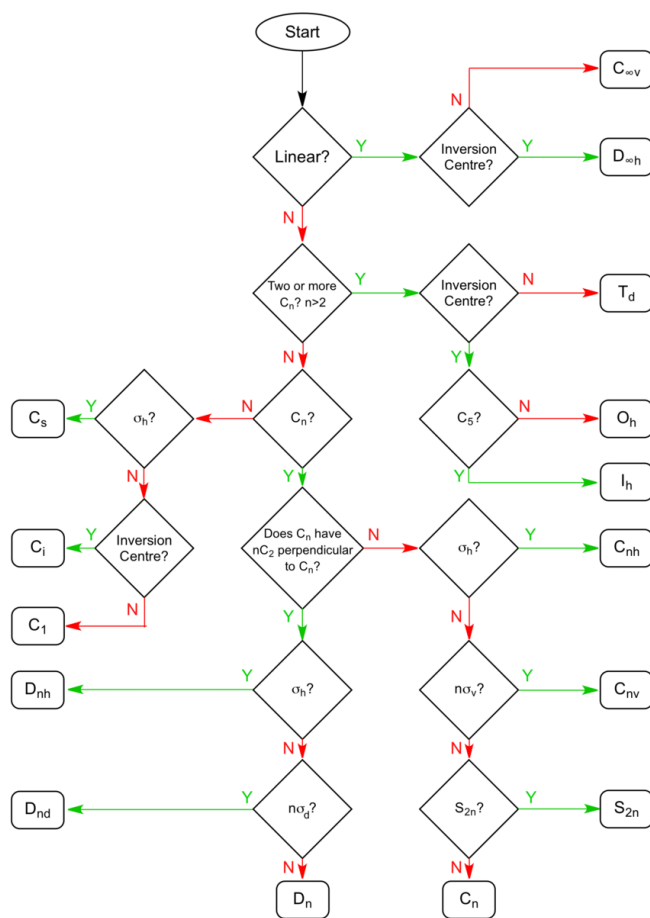


Figure 10: Decision tree for determining the point group of a molecule

## 2 Observations

**Observation 1.** *Machine Learning (ML) and Artificial Intelligence (AI) are not the same concept, in fact, ML models are a subset of AIs. The key relies in word Learning, generally speaking, an AI does not need to actually learn from a set of data. It can be set up within a decision tree such that logically responds with the proper answer following its criteria.*

*However a Machine Learning Method is a kind of AI that learns from data and evolves with the provided data.*

**Observation 2.** *Mathematically speaking, Machine Learning methods are no different from an optimisation problem.<sup>13</sup>*

*Generally speaking, all ML methods seek to find a criterion to determine the quantity of a certain property of some data, which we lack knowledge of and will call the target. In order to do so, they rely on a dataset where the sought property is already known and which establishes relationships between this property and other features of the data. Then, based on these relationships and the properties of the target, the model searches for the position where the target fits within the data and determines the quantity of the property, taking this placement into account within the known data.*

*For instance, we could think of standard calibration with linear regression as an ultra-simplification of a Machine Learning problem. We have a set of data (such as salt concentration in a solution)  $\{x_i\}_{i=0}^n$ , we also know the quantity of a certain property (such as the solution's conductivity) for each entry  $x_j$ , denoted as  $\{f(x_i)\}_{i=0}^n$ . In this example, our ML method assumes that the data and the property are linearly related and that the distance between data and the property is equivalently defined as  $d(x_i, x_j) = \sqrt{x_i^2 + x_j^2}$ . Hence, based on the assumption of linearity, the best possible relationship is a straight line, and the method*

*looks for the line that minimises the distance between the line and the data.*

*Then, given a target (for instance, a certain concentration), about which we initially lack knowledge, the method is capable of quantitatively computing the theoretical property (for instance, its conductivity) by using the fitted line representing the relationship. Consequently, this can be helpful in determining the theoretical concentration of salt needed to achieve a certain conductivity without having actual data for that concentration.*

*The learning part (the most important part) comes from the fact that the model (the line) learns from the data. As we train our system (i.e., perform linear regression) with increasingly larger datasets, the accuracy of our model's predictions improves. Thus, we can conclude that the system is learning from the data, which qualifies it as a Machine Learning model.*

**Observation 3.** *Since we are working with molecules, our Machine Learning model will require certain parameters to determine the relationship with the property.*

*In this case, our property (i.e.,  $\{f(x_i)\}_{i=0}^n$ ) will be the inhibition of a specific protein (COX-2), and more specifically, it will be  $IC_{50}$ ,  $IC_{90}$ , or  $IC_{99}$  (1).*

*Moreover, our parameters (i.e.,  $\{x_i\}_{i=0}^n$ ) will be the molecular descriptors, which must be **carefully chosen**. This problem introduces additional complexities compared to the linear regression case:*

- *We will be working with approximately  $10^3$  descriptors, making this a  $10^3$ -dimensional problem.*
- *We want the system to automatically discard non-relevant data, as some molecules could be exceptions due to external factors.*

- We do not assume linearity.
- Not all descriptors have the same distances or weights, as we will want some of them to be more relevant than others.
- Some descriptors will be continuous (for instance, bond distances) while others will remain discrete (for instance, topological indexes) by construction which implies that some dimensions need to be treated differently to others.

Interestingly enough, despite these complexities, in some way, we can say that this project is not entirely different from computing a linear regression...

**Observation 4.** Taking into account definitions (6) and (7), we can also define a norm, denoted by  $\mathcal{N}(A, B)$ , which refers to the norm between molecules  $A$  and  $B$ . For instance, this can be expressed as  $\mathcal{N}(A, B) = |\text{Ind}_A - \text{Ind}_B|$ , where  $\text{Ind}_A$  and  $\text{Ind}_B$  represent the respective indices (Wiener, Zagreb, etc.) of the molecules<sup>5</sup>. Using this norm, we can generate a metric space that tells us how close two molecules are, based on the similarity of their indices.

In this way, we can apply the concept of “linear regression” and compute the relationship between the indices and the  $\text{IC}_{50}$ . This refers to one dimension, corresponding to a specific topological index. We can then extend this to consider all remaining indices (approximately  $10^3$ ), and find the optimal position in this  $10^3$ -dimensional space for our resulting molecule.

**Observation 5.** We have extracted all the available data from the ChEMBL database with a known  $\text{IC}_{50}$  value for the COX-2. After deleting entries with the same canonical smile 4827 entries remain. By keeping just those with a

$\text{IC}_{50} \in [-1, 200]^6$ , we obtain a remaining set of 1518 entries represented in the Figure (11).

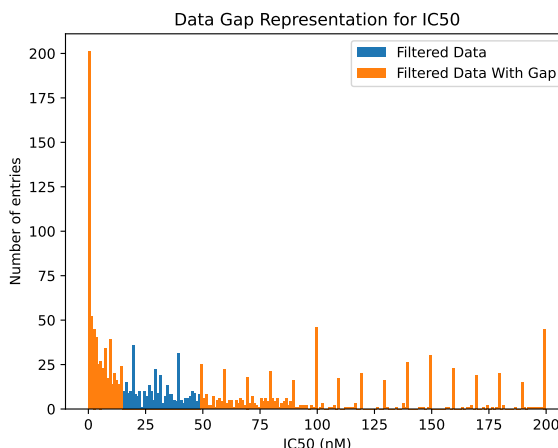


Figure 11: Data representation after the filtering.

**Observation 6.** Figure (11) clearly shows multiples of 10 are significantly overrepresented respect to other values. Additionally, by analyzing data from other proteins we can confirm this is not an exception but rather the norm. It looks like experimental chemists have a preference for multiples of 10, often rounding their results to the nearest value in this category. Consequently, an inherent error associated with these values will always be present, and we will have to deal with it.

**Observation 7.** Once the chemical descriptors are computed, we proceed to calculate the Pearson correlation coefficients as defined in Definition (42). The results are subsequently plotted in Figure (12). From the plot, we observe a “Gaussian bell” shape centered around zero, which aligns with our expectations. This suggests that there is no inherently preferred set of chemical descriptors for the protein, a finding that may initially seem at odds with the

<sup>5</sup>Actually, any kind of norm would work, as long as it depends on the indices.

<sup>6</sup>Note that the lower constrain won’t be ever surpassed by the  $\text{IC}_{50}$  definition 1. This is equivalent to avoid constraining the data.

project's goals.

However, our hypothesis states that there exists a combination (or combinations) of chemical descriptors that are directly correlated with the inhibition of the protein. This observation does not contradict that statement. Intuitively, we believe that filtering out descriptors with the lowest  $r^2$  values might lead to improved results. This possibility will be explored in subsequent discussions.

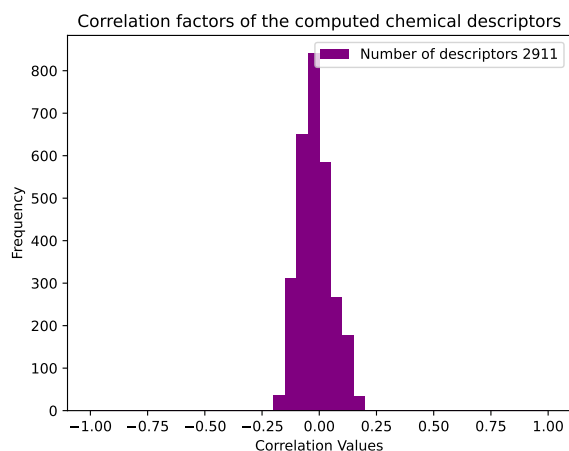


Figure 12: Representation of the Pearson correlation factors.

**Observation 8.** The first algorithm implemented in this project is the Random Forest machine learning algorithm. Although the project's primary focus is not the detailed workings of machine learning algorithms, it is natural to ponder the essence of what the machine is doing. The name "Random Forest" is directly tied to the analogy between a tree and a decision tree.

As described in Definition (43), the Random Forest algorithm generates a collection of random decision trees. The number of trees in the forest is a configurable parameter of the algorithm. More specifically, the Random Forest algorithm produces a set of decision trees, each trained on a random subset of the training data.

The structure of each tree adjusts during training to make predictions that closely approximate the true target values. The final output of the Random Forest is typically an aggregate of the predictions from all the individual trees, such as a majority vote for classification tasks or an average for regression tasks.

One notable characteristic of this method is its inherent randomness. Each Random Forest model is unique, even if configured with the same parameters and initialized under identical conditions. However, experimental evidence suggests that this randomness does not significantly impact the results. If two Random Forest algorithms are trained on the same dataset using identical parameters, their outputs will converge, even if the internal structures of their decision trees differ.

**Observation 9.** From Observation (7), one might hypothesize that training with chemical descriptors that have higher correlation factors, or discarding those with low correlation factors, would lead to better predictions. Numerically, however, this hypothesis is contradicted, as shown in Table (1). Despite this, the result remains consistent with our core hypothesis: the correlation factor of a chemical descriptor does not necessarily indicate its relevance in protein inhibition. Our key assertion is that there exists a combination (or combinations) of chemical descriptors that are directly correlated with the inhibition of the protein.

This suggests that the specific combinations of descriptors correlated with protein inhibition are not dependent on the individual correlation factor values. At first glance, a lack of correlation does not imply a lack of relationship to the inhibition!

**Observation 10.** The dependency on the gap generations seems much relevant than the dependency on the descriptors' correlation factor (Observation (9)). By changing the erased percentage (i.e. the fraction of data discarded in

the gap generation) we can compute  $r^2$  and the Mean Square Error (i.e. the statistical deviation) which provides the results depicted in Figure (13) and Figure (14) respectively.

Figure 13: Plot of  $\%_{\text{Erased}}$  vs  $r^2$ . For: maximum  $IC_{50} = 200$  nM, taking into account all the chemical descriptors, the 20.0% of the remaining data is saved for testing and the 80.0% for training. The training is done with 200 trees.

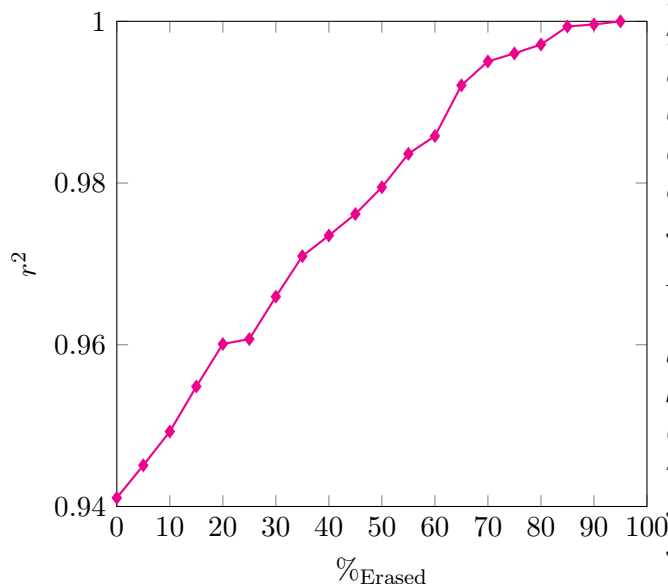
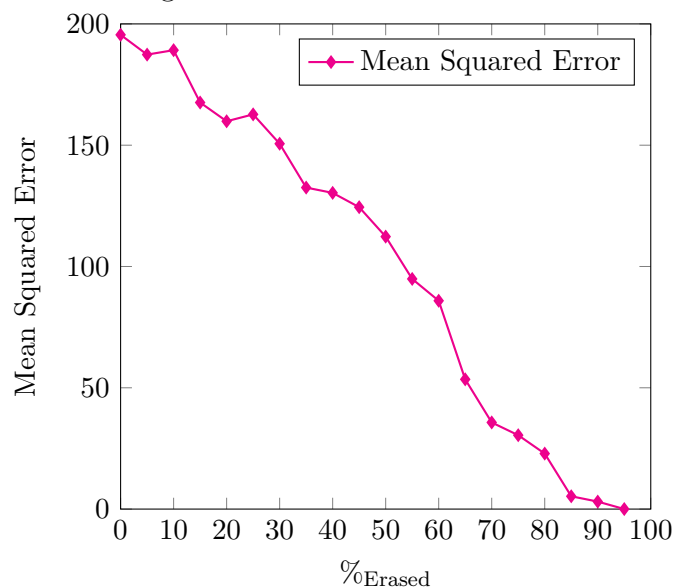


Figure 14: Plot of  $\%_{\text{Erased}}$  vs Mean Squared Error. For: maximum  $IC_{50} = 200$  nM, taking into account all the chemical descriptors, the 20.0% of the remaining data is saved for testing and the 80.0% for training. The training is done with 200 trees.



**Observation 11.** From Observation (6), we know the data has an inherent precision problem, especially when approaching higher  $IC_{50}$  values, additionally we do not know the precision of some measurements. Since our working models do not account for the lack of precision in each molecule's data, the values of  $r^2$  and mean squared error are not expected to be exactly 1 and 0, respectively. In fact, some variance and uncertainty are inevitable. Therefore, discarding more than 50% of our data would be a critical mistake, even if statistical measures suggest otherwise.

For example, suppose we aim to build a table. We have an arbitrarily large number of table legs, each varying slightly in length. On one hand, if we choose to build a three-legged table, we would obtain a rigid structure with no variance (i.e. statistically, this would be considered the "best" table). However, it would be steeply inclined and therefore useless, as any object placed on it would slide off. On the other

hand, if we use as many legs as possible, the table may initially be unstable due to variations in leg length. Nevertheless, as more legs are added, the instability decreases, eventually allowing us to place objects on it and even have a meal.

This example highlights the fact that an  $r^2$  value of 1 and a mean squared error of 0 are not absolute requirements in this case, due to Observation (6). However, these values still serve as a good reflection of our method's reliability up to a certain threshold.

In this work, we will arbitrarily set this threshold at  $r^2 = 0.96$ . This hard-coded decision is purely conventional, as we need to establish a point beyond which  $r^2$  values become meaningless (i.e., obtaining  $r^2$  values above 0.96 is unnecessary, given that such variance is inherent to our data). The justification for this choice is that, in some areas of chemistry<sup>7</sup> two data sets are generally considered correlated if their  $r^2$  value exceeds 0.96.

However, it is important to emphasize that this convention is purely subjective and should not be regarded as an inviolable rule.

---

<sup>7</sup>Particularly in analytical chemistry



### 3 Reflections

**Reflection 1.** *During the first meeting with Dra. González, Dr. Moreno & Dr. Lluch a certain theme was brought up. "Wouldn't be great to actually have the wave function of the molecules?" This way we could avoid the use of descriptors and QSARs and have more accurate results.*

*Though it was a joke, it might not be a utopia someday. I believe quantum computers will be capable of this task and maybe a hundred years*

*from today, this project could be doable working directly with the waves functions (not analytically obviously). How knows? Time will tell us.*

**Reflection 2.** *Fun fact,<sup>3</sup> is a mathematical article, completely apart from the chemical world but somehow we can manage to use it to describe molecules whatsoever. I believe applying knowledge from other sectors is key to evolve and reach even greater results.*

### 4 Packages

**Package 1.** *"Requests allows you to send HTTP/1.1 requests extremely easily. There's no need to manually add query strings to your URLs, or to form-encode your POST data. Keep-alive and HTTP connection pooling are 100% automatic, thanks to urllib3."*<sup>14</sup> *We will use this package to create requests to the ChEMBL database.*

**Package 2.** *"Pandas is a Python package that provides fast, flexible, and expressive data structures designed to make working with "re-*

*lational" or "labeled" data both easy and intuitive. It aims to be the fundamental high-level building block for doing practical, real world data analysis in Python. Additionally, it has the broader goal of becoming the most powerful and flexible open source data analysis / manipulation tool available in any language. It is already well on its way towards this goal"*<sup>15</sup> *We will use this package to interact with the data using a built-in Pandas class called Data Frames.*

### 5 Physical books in the library

**Book 1.** *Molecular descriptors for chemoinformatics. Todeschini, Roberto.; Consonni, Viviana. 2009. Volum: 2.*

*Localization: Ciència i Tecnologia.*

*Signatura: 54:68 Tod.*

*Codi de Barres: 1501208901*

## 6 Results

Table 1: Results for; Maximum  $IC_{50} = 200$  nM, Gap of the 25.0% of the set, the 20.0% of the remaining data is saved for testing and the 80.0% for training. The training is done with 200 trees.

|                       | $\sigma \leq 0$    | $\sigma \leq 0.05$ | $\sigma \leq 0.075$ | $\sigma \leq 0.10$ |
|-----------------------|--------------------|--------------------|---------------------|--------------------|
| $r^2$                 | 0.9607020297015885 | 0.9608545198149474 | 0.9605224652065423  | 0.9599158186858139 |
| Mean Squared Error    | 162.70220420738363 | 162.0708617391865  | 163.4456405456867   | 165.9572900112382  |
| number of descriptors | 2911               | 1485               | 1001                | 568                |

Table 2: Results for; Maximum  $IC_{50} = 200$  nM, taking into account all the chemical descriptors, the 20.0% of the remaining data is saved for testing and the 80.0% for training. The training is done with 200 trees.

| %Erased            | 0%                 | 5%                 | 10%                | 15%                  |
|--------------------|--------------------|--------------------|--------------------|----------------------|
| $r^2$              | 0.9410691012016194 | 0.9450975181635125 | 0.9492662637374885 | 0.954840165823525    |
| Mean Squared Error | 195.50415499282306 | 187.34588735050227 | 189.17427179281864 | 167.5904134577239    |
| %Erased            | 20%                | 25%                | 30%                | 35%                  |
| $r^2$              | 0.960093147530398  | 0.9607020297015885 | 0.9659474255238565 | 0.970950038525679    |
| Mean Squared Error | 159.86610568752263 | 162.70220420738363 | 150.60605102824292 | 132.53106869928666   |
| %Erased            | 40%                | 45%                | 50%                | 55%                  |
| $r^2$              | 0.9735150867866535 | 0.9761682347024537 | 0.9794766777874534 | 0.9836139552920377   |
| Mean Squared Error | 130.3450775382814  | 124.45078192954777 | 112.31398471818515 | 94.8719409443141     |
| %Erased            | 60%                | 65%                | 70%                | 75%                  |
| $r^2$              | 0.9858128266667714 | 0.9920845685273683 | 0.9950402754045806 | 0.9960292052079258   |
| Mean Squared Error | 85.87640609411348  | 53.48552498895907  | 35.71100434005104  | 30.45304023300008    |
| %Erased            | 80%                | 85%                | 90%                | 95%                  |
| $r^2$              | 0.9971354719623441 | 0.9993589449322425 | 0.9996146336878613 | 0.9999982353685487   |
| Mean Squared Error | 22.873087477997174 | 5.275980426255672  | 3.083365552988717  | 0.010325883398321972 |

## 7 Various

Listing 1: Example of a JSON file.

```
1 {
2   "page_meta": {
3     "limit": 1000,
4     "offset": 0,
5     "total_count": 1500,
6     "next": "/chembl/api/data/activity.json?limit=1000&offset=1000&
7     target_chembl_id=ChEMBL372&standard_type=IC50",
8     "previous": null
9   },
10  "activities": [
11    {
12      "activity_id": 123456789,
13      "assay_chembl_id": "ChEMBL123456",
14      "molecule_chembl_id": "ChEMBL25",
15      "canonical_smiles": "CCOCCO",
16      "standard_value": "50",
17      "standard_units": "nM",
18      "standard_type": "IC50",
19      "target_chembl_id": "ChEMBL372",
20      "assay_type": "B",
21      "relation": "=",
22      "document_chembl_id": "ChEMBL12345"
23    },
24    {
25      "activity_id": 987654321,
26      "assay_chembl_id": "ChEMBL654321",
27      "molecule_chembl_id": "ChEMBL1095",
28      "canonical_smiles": "CCCN(CC)CC",
29      "standard_value": "25",
30      "standard_units": "nM",
31      "standard_type": "IC50",
32      "target_chembl_id": "ChEMBL372",
33      "assay_type": "B",
34      "relation": "<",
35      "document_chembl_id": "ChEMBL54321"
36    }
37  ]
38 }
```

## References

- [1] Swinney, D. C. In *Chapter 18 - Molecular Mechanism of Action (MMoA) in Drug Discovery*; Macor, J. E., Ed.; Annual Reports in Medicinal Chemistry; Academic Press, 2011; Vol. 46; pp 301–317.
- [2] Gasteiger, J.; Engel, T. *Molecular Descriptors for Chemoinformatics*; Wiley-VCH: Weinheim, Germany, 2003.
- [3] Xu, K. *Applied Mathematics Letters* **2011**, *24*, 1026–1030.
- [4] Randic, M. *Journal of the American Chemical Society* **1975**, *97*, 6609–6615.
- [5] Gutman, I.; Furtula, B.; Katanić, V. *AKCE International Journal of Graphs and Combinatorics* **2018**, *15*, 307–312.
- [6] Todeschini, R.; Consonni, V. *Molecular Descriptors for Chemoinformatics*, 2nd ed.; Wiley-VCH: Weinheim, Germany, 2009.
- [7] Burden, F. R.; Wilkins, D. C. *Journal of Chemical Information and Computer Sciences* **1996**, *36*, 1–5.
- [8] Güner, O. F. *Pharmacophore Perception, Development, and Use in Drug Design*; International University Line, 2000.
- [9] National Cancer Institute Definition of COX-2 - NCI Dictionary of Cancer Terms. 2024; <https://www.cancer.gov/publications/dictionaries/cancer-terms/def/cox-2>, Accessed: 2024-10-09.
- [10] Commons, W. File:6COX.png — Wikimedia Commons, the free media repository. 2023; <https://commons.wikimedia.org/w/index.php?title=File:6COX.png&oldid=770187485>, [Online; accessed 9-October-2024].
- [11] Davies, N. M.; Jamali, F. *Pharmacology & Therapeutics* **2000**, *89*, 133–155.
- [12] Commons, W. COX-2 receptor site and its amino acid profile with celecoxib in the binding site. 2024; <https://commons.wikimedia.org/w/index.php?title=File:COXII-receptor.png&oldid=844592069>, [Online; accessed 12-October-2024].
- [13] Deisenroth, M. P.; Faisal, A. A.; Ong, C. S. *Mathematics for Machine Learning*; Cambridge University Press, 2020.
- [14] Reitz, K.; Chalasani, A. Requests: HTTP for Humans. <https://pypi.org/project/requests/>, 2023; Python package.
- [15] McKinney, W. pandas: A Foundational Python Library for Data Analysis. 2023; <https://pandas.pydata.org>, Version 1.5.3.