# A QSAR STUDY ON THE MUTAGENIC ACTIVITY OF AZO DYES AND AROMATIC AMINE DERIVATIVES

by

Müge Küçükali

BS. in Environmental Engineering, Istanbul Technical University, 2013

Submitted to the Institute of Environmental Sciences

in partial fulfillment of the requirements for the degree of

Master of Science

in

Environmental Sciences

Boğaziçi University

2019

# A QSAR STUDY ON THE MUTAGENIC ACTIVITY OF AZO DYES AND AROMATIC AMINE DERIVATIVES

APPROVED BY:

Prof. Dr. Melek Türker Saçan . . . . . . . . . . . . . .
Thesis Advisor

Prof. Dr. Nadim Copty . . . . . . . . . . . . . .

Prof. Dr. Safiye Sağ Erdem . . . . . . . . . . . . . .

DATE OF APPROVAL: 31/10/2019

# ACKNOWLEDGEMENTS

# ABSTRACT

# A QSAR STUDY ON THE MUTAGENIC ACTIVITY OF AZO DYES AND AROMATIC AMINE DERIVATIVES

In the present study, the mutagenic activity of diverse structure of azo dyes and aromatic amine derivatives on TA98 *Salmonella typhimurium* bacterial strain with S9 activation was used to generate a quantitative structure – activity relationship (QSAR) model. The descriptors required for the model development were obtained by SPARTAN (v.10), DRAGON (v.7.0) software packages. The selection of descriptors was carried out by the tools implemented in QSARINS (v.2.2.3) software. Many division trials were performed on the dataset as training and test sets which comprise the 80% and 20% of the whole set, respectively. 6 descriptors (CIC2, Chi_D/Dt, L/Bw, TDB09p, Mor28s and piPC08) constitute the final model. The applicability domain (AD) of the generated QSAR model was defined by both the ranges of response and descriptors. The predictive ability of the final model was tested using an external dataset consisting of currently used 33 anionic water soluble textile dyes (eight anthraquinone and 25 azo dyes) with no experimental mutagenicity data. The proposed QSAR model had 70% structural coverage for the external set chemicals. The mutagenic activities of 22 current textile dyes belonging anthraquinone and azo dyes fell within the applicability domain of the proposed QSAR model which means that they were well predicted by the model. The order of 6 anthraquinone dyes which are mainly used for cotton, fiber dyeing and leather shading with the most mutagenic activity is: Acid Blue 62>Acid Blue 40>Acid Blue 45>Acid Blue 80>Acid Blue 230>Acid Blue 344. The least mutagenic azo dye is Direct Orange 34 which is mainly used for cotton, silk, wool and their blended fabric dyeing and printing, also can be used for leather and paper shading.

# ÖZET

## AZO BOYALARI VE AROMATİK AMİN TÜREVLERİNİN MUTAJENİK AKTİVİTESİ ÜZERİNE BİR KYAİ ÇALIŞMASI

Bu çalışmada, çeşitli azo boyaları ve aromatik amin türevlerinin S9 aktivasyonlu TA98 *Salmonella typhimurium* bakteri türü üzerindeki mutajenik aktivitesi kullanılarak kantitatif yapı - aktivite ilişkisi (KYAİ; (QSAR)) modeli geliştirilmiştir. Model geliştirme için gerekli tanımlayıcılar SPARTAN (v.10) ve DRAGON (v.7.0) yazılım programları ile elde edilmiştir. Tanımlayıcıların seçimi, QSARINS (v.2.2.3) yazılımında bulunan araçlar ile yapılmıştır. Tüm veri setinin % 80'i ve % 20'si sırası ile eğitim ve test setleri olacak şekilde farklı ayrımlar denenmiştir. 6 (CIC2, Chi_D/Dt, L/Bw, TDB09p, Mor28s and piPC08) tanımlayıcı ile nihai modeli oluşturulmuştur. Nihai modelin tahmin yeteneği, deneysel mutajenite verisi olmayan, halihazırda kullanılan 33 anyonik suda çözünür tekstil boyalarını (sekiz antrakinon ve 25 azo boyası) içeren harici bir veri seti kullanılarak test edilmiştr. Oluşturulan modelin uygulanabilirlik alanı (AD), hem aktivite hem de tanımlayıcıların oluşturduğu aralık gözönünde bulundurularak tanımlanmıştır. Önerilen KYAİ modelinin yapısal olarak dış setteki bileşiklerin %70 'ini kapsadığı görülmüştür. Halen kullanımda olan 22 adet antrakinon ve azo boyasının mutajenik aktiviteleri, önerilen QSAR modelinin uygulanabilirlik alanı içinde kalmıştır. Bu da mutajenik aktivitelerinin model tarafından iyi tahmin edildiklerini göstermiştir. En yüksek mutajenik aktivite gösteren, pamuk, elyaf boyama ve deri gölgemesi için kullanılmakta olan 6 antrakinon boyanın mutajenik aktivitelerinin sıralaması Asit Mavi 62>Asit Mavi 40>Asit Mavi 45>Asit Mavi 80>Asit Mavi 230>Asit Mavi 344 'tür. En düşük mutajenik aktiviteye sahip azo boya ise esas olarak pamuk, ipek, yün ve bunların karışımlı kumaş boyaması ve baskısı için kullanılmakta olan Direkt Turuncu 34' tür. Bu boya ayrıca deri ve kağıt gölgelemesi için de kullanılmaktadır.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF SYMBOLS/ABBREVIATIONS

| Symbol | Explanation | Unit |
|---|---|---|
| μ | Dipole Moment | |
| η | Hardness | |
| ω | Electrophilicity index | |
| $\sigma$ | Softness | |
| $E$ | Gas-phase Energy | eV |
| $E_{aq}$ | Aqueous-phase Energy | eV |
| $E_{HOMO}$ | Energy of the Highest Occupied Molecular Orbital | eV |
| $E_{LUMO}$ | Energy of the Lowest Unoccupied Molecular Orbital | eV |
| $R^2$ | Coefficient of Determination | |
| $R^2_{adj}$ | Adjusted (for degrees of freedom) $R^2$ | |
| Chi_D/Dt | Randic-like Index from Distance/Detour Matrix | |
| L/Bw | Length-to-breadth Ratio by WHIM | |
| CIC2 | Complementary Information Content Index (neighborhood symmetry of 2-order) | |
| Mor28s | Signal 28 / weighted by I-state | |
| TDB09p | 3D Topological distance based descriptors lag 9 weighted by polarizability | |
| piPC08 | Molecular Multiple Path Count of Order 8 | |
| WHIM | Weighted Holistic Invariant Molecular Descriptors | |
| $h*$ | Critical Hat Value | |

| Abbreviation | Explanation |
|---|---|
| AD | Applicability Domain |
| CAS | Chemical Abstracts Service |
| $CCC$ | Concordance Correlation Coefficient |
| EC | European Commission |

| | |
|---|---|
| ECHA | European Chemicals Agency |
| EU | European Union |
| $F$ | Fischer Statistics |
| Log$P$ | Logarithm of $n$-Octanol/Water Partition Coefficient |
| $MAE$ | Mean Absolute Error |
| MCDM | Multiple Criteria Decision Making |
| MLR | Multiple Linear Regression |
| OECD | Organization of Economic Co-operation Development |
| PM6 | Parameterized Model 6 |
| $PRESS$ | Predicted Residual Sum of Squares |
| QSAR | Quantitative Structure-Activity/Toxicity Relationship |
| $Q^2_{\text{LOO}}$ | Leave-One-Out Cross Validation |
| $Q^2_{\text{LMO}}$ | Leave-Many-Out Cross Validation |
| $QUIK$ | $Q^2$ Under Influence of $K$ |
| $RMSE$ | Root Mean Square Error |

# 1. INTRODUCTION

Dyes have been applied in many industries for a long time, including the textile, printing, medical and energy industries. There are some properties of dyes which can be tested experimentally such as colour, brightness, solubility, fastness, mutagenicity, diffusion constant and so on (Kothari et al., 2013). Some of the dyes are toxic, carcinogenic, mutagenic or harmful to human health. The textile industry accounts for the largest consumption of dyestuffs at nearly 70 percent (Mathur et al., 2012). Textile dyes are mainly classified according to their uses in the dyeing process. Many functional classes of dyes, containing acid, basic, disperse, reactive and solvent dyes comprise azo compounds (Freeman, 2013). Azo dyes account for the largest proportion of all synthetic dyes and include approx. 70 % of all organic dyes, which are currently on the market and are manufactured mainly in China, India, Korea, Taiwan, and Argentina. Clothing textiles on the European market were mostly dyed in those Asian countries, and are then imported to Europe.

Aminoazo derivatives are important because of their widespread use in the textile industry (Stead, 1990). Some azo dyes are both mutagenic and carcinogenic. It has been proved that variety of 4-aminoazobenzene (AAB), *N*-methyl-4-aminoazobenzene (MAB) and *N*,*N*-dimethyl-4-aminoazobenzene (DAB) derivatives are mutagenic (Garg et al., 2002). It was known that benzidine is a mutagenic moiety of many azo dyes and can be generated from azo dyes through the reduction by intestinal and environmental microorganisms (Chung et al., 2006). Azo dyes also produce free aromatic amines that are significantly mutagenic and carcinogenic. In addition to the effects caused by exposure to contaminated water and food, workers who deal with these dyes can be exposed to them. Moreover, scientific researches have proven their adverse health effects. Myslak et al. (1991) observed that German painters developed bladder cancer after long- time exposure to azo dyes.

Registration, Evaluation, Authorization and Restriction of Chemicals (REACH) is the Regulation (EC) with the No 1907/2006 of the European Parliament and of the Council of 18 December 2006. It came into force on 1 June 2007. The purpose of this regulation is to protect human health and the environment. The enforcement of REACH has a great impact on the textile and leather industry.

Textile and leather products belong to articles under REACH regulation. Those products might contain Substances of Very High Concern (SVHC) and REACH restricted substances regulated by

REACH regulation. Suppliers (manufacturers, importers or retailers) shall first identify all possible toxic chemicals in their products. If the products do not contain SVHC and REACH restricted substances above the threshold limits of REACH, the products are compliant with REACH.

If products contain any SVHC and restricted substances above the threshold limits, suppliers are responsible to identify these hazardous substances belong to (SVHC candidate list or REACH restriction list or Authorization list) and adopt the required measures to achieve REACH compliance.

Of the minimum required data for the assessment of environmental and human hazards under REACH regulation CMR (Carcinogenic, Mutagenic and Reproductive toxicity) properties are of major concern together with PBT (Persistence, Bioaccumulation and ecoToxicology) properties. However, this information is not available for the majority of the existing chemicals.

Mutagenicity is widely recognized as a valid surrogate of carcinogenicity. The term mutagenicity refers to the ability of a chemical to induce genetic damages that may occur by several mechanisms involving interactions with the DNA or with both the DNA and other cellular targets, e.g. proteins (i.e., chromosomal aberrations, and changes in the number of chromosomes) (Benigni et al., 2011).

Mutagenic activity is measured with Ames test which is known as the bacterial reverse mutation assay (Ames et al., 1973a; Ames et al., 1973b). It is a well-known and widely used assay to detect mutagenicity *in vitro* and is of crucial importance in development as an early alerting system for potential carcinogenicity and/or teratogenicity. Kaur (1993a and 1993b) screened textile azo dyes for mutagenicity with Ames/Salmonella assay with and without metabolic activation. Many of the amines were found to be mutagenic in tester strains TA98 and TA100 but require exogenous mammalian activation (S9) for activity (Chung, 2006).

The information concerning the mechanism of action of the molecules studied thanks to the various strains of *Salmonella typhimurium* available. In the Ames mutagenicity, histidine dependent strains of *Salmonella typhimurium* is used. When these strains are exposed to a mutagen, reverse mutations that restore the functional capability of the bacteria to synthesize histidine enable bacterial colony growth on a medium deficient in histidine ("revertants"). Knowing that many chemicals interact with genetic material only after metabolic activation by enzyme systems not available in the bacterial cell, the test compounds are in many cases additionally examined in the presence of a mammalian metabolizing system, which contains liver microsomes (rat liver S9 mix). TA98 bacterial strain with S9 activation is well known to detect frameshift mutagens (Carroll et al., 2002). Some of

azo dyes such as dimethylamino-phenylazoenthiazole (6BT) are mutagen in the *Salmonella typhimurium* TA98 bacterial tester strain in the presence of an induced rat-liver S9 mix (TA98+S9) (Sztandera, 2003). Kojima et al. (1991) showed that the 3 methoxy-4 aminoazo-benzene (3-OMe-AAB) is a strong mutagen in *Escherichia coli* and *Salmonella typhimurium*, whereas 2-OMe-AAB is an extremely weak mutagen under the same condition.

The mutagenic potency of azo dyes is strongly dependent on the position of substituents with respect to both the aromatic rings and the amino nitrogen atom. This implies that there is a relationship between chemical structures and mutagenic activities. Mutagenicity data for most of the chemicals in use including dyes are scarce and should be known to understand their effects on the environment. However, conducting experiments for this purpose is time and energy consuming in addition to being costly. For this reason, quantitative structure-activity relationship (QSAR) method for predicting mutagenic activities without any experimental data may be desirable. Due to the lack of experimental data, there has been increasing use of theoretical structural descriptors in the hazard assessment of chemicals via QSAR models. QSAR studies are based on the idea that property/biological activity of a chemical can be predicted or characterized by its structure. Predicted property/biological activity values from valid QSAR models have a high potential to fill the data gaps. Another major goal of QSAR is the rational design of the new and more active compound starting from the known ones. One advantage is that this method needs no equipment and reagent. Researchers compute the molecular descriptors by computer using chemical structures and software packages. QSAR models have been used for decades and the usage of QSAR is recommended and supported by European Chemical Agency (ECHA, 2017).

In the literature, although there are some QSAR studies on the prediction of mutagenic activity of dyes (Garg et al., 2002; Bhat et al., 2005), they have some drawbacks. As such, Pasha et al. (2008) developed several QSAR models for the mutagenicity of aminoazobenzenes which do not carry out the principles set by Organization for Economic Co–operation and Development (OECD, 2007) and not tested with the up-to-date validation criteria. To increase the regulatory acceptance of (Q)SAR methods, the OECD is developing a QSAR Toolbox to make (Q)SAR technology readily accessible, transparent, and less demanding in terms of infrastructure costs in 2008. The Toolbox incorporates information and tools from various sources into a logical workflow. Devillers et al. (2010) evaluated the toolbox for estimating the mutagenicity of chemicals. Moreover, VEGA which is another ideal application for batch processing large dataset provides freely tens of QSAR models to predict the properties such as persistence, logP, bioconcentration factor (BCF), carcinogenicity, mutagenicity,

skin sensitization of the compounds. Instituto Mario Negri is the coordinator of the VEGA project (www.vegahub.eu, last accessed June 2019).

Mombelli et al. (2016) studied the mutagenicity of the compounds using the different types of QSAR models such as VEGA QSAR. This application can be used in any system supported by JAVA. As it is mentioned above, the application of QSAR in filling the data gap is supported by the European Chemical Agency (ECHA). ECHA has reorganized its Practical Guide on How to use and report QSARs with further advice and examples on using QSARs for registering under REACH in Helsinki, 17 March 2016.

While the toxic effects of many of the chemicals now entering the public market place are being routinely examined, many thousands of new and untested compounds are being synthesized annually in thousands of laboratories and scientists are becoming increasingly concerned about the potential toxicity of the new chemicals they are being exposed to. Therefore, scientists focus on the design and development of non-toxic and environmentally friendly consumer goods using the principles of green chemistry (Clark, 2006). Hence, QSAR also promotes green chemistry. Not only due to their higher efficiency and lower cost, but they can also obtain a rapid assessment of the potential effects of chemicals on human health and also the environment. Therefore, the QSAR models have been developed as feedback to different legislation around the world (e.g. REACH) to help in decreasing animal testing and designing more greener compounds. Many reliable models have been produced by this method, and it was also employed in the present study. In an attempt to find an adaptable, well-performing and predictive model for mutagenicity of dyes, we applied QSAR modelling to predict the unknown mutagenic activity of chemicals (particularly current textile dyes) using the mutagenic activity of aminoazo dyes and aromatic amine derivatives. Such studies are useful in terms of identification and prioritization of dangerous chemicals as well as providing data for a wide range of chemicals with no experimental mutagenicity data.

## 1.1. Purpose of the Study

The aim of the present study is; (1) to compile mutagenicity data for azo dyes and aromatic amine derivatives from literature, (2) to develop valid linear QSAR models for the mutagenic activity for aminoazodyes and aromatic amine derivatives in the TA98 *Salmonella typhimurium* bacterial strain with S9 activation, (3) to evaluate the performance of the developed QSAR models using a test/validation set, in which data were not used in the generation of the model, and (4) to apply all the up-to-date internal and external validation criteria to the generated QSAR models, and (5) to test the

predictive performance of the selected  QSAR model for the unknown mutagenic activity of currently used 33 anionic water-soluble textile dyes (eight anthraquinone and 25 azo dyes).

# 2. THEORETICAL BACKGROUND

## 2.1. Dyes

Dyes are classified according to their application and chemical structure and are composed of a group of atoms known as chromophores, responsible for the dye color. These chromophore-containing centers are based on diverse functional groups, such as azo, anthraquinone, nitro, aryl-, carbonyl and others (Prasad, 2010). With respect to the number and production volumes, azo dyes are the largest group of colorants, constituting 60-70% of all organic dyes produced in the world (Bafana, 2011). An early theory of dyes first formulated by O. Witt provided a basis for understanding the relation between color and the structure of the compound. According to this theory, a dye consists of three components: one or more the fused benzene rings attached to the unsaturated groups called as chromophores (e.g., $-N=N-$, $-NO_2$, $-C=O$) and basic groups called as auxochromes (e.g., $NH_2$, $OH$ groups). A chromophore is a chemical group that is responsible for the appearance of color in compounds where it is located. The colorants are sometimes also classified according to their main chromophore (e.g., azo dyes contain the chromophore $-N=N-$) (Iqbal, 2008).

Although dyes have an affinity to the substrates to which they are being applied (Pereira and Alves, 2012), they colored substances which are soluble or go into solution during the application process (Bafana et al., 2011). As such, up to 200000 tons of these dyes are lost to effluents every year during the dyeing and finishing operations, due to the inefficiency of the dyeing process (Maas and Chaudhari, 2005; Jin et al., 2007; Ogugbue et al., 2011; Saratale et al., 2011). Unfortunately, most of these dyes escape conventional wastewater treatment processes and are often found chemically unchanged in wastewater treatment plant effluents; contaminating raw water or become concentrated in the sludge, causing a disposal problem (Carneiro et al., 2010). Very small amounts of dyes in the water to cause a highly visible change in color. They can also affect the aquatic ecosystem, decreasing the passage of light penetration and gas dissolution in lakes, rivers and other bodies of water (Saranaik and Kanekar, 1995; Banat et al., 1996; Modi et al., 2010). Shaul et al. (1991) studied 18 azo dyes and found that 11 passed practically unchanged through the activated sludge system, 4 were adsorbed by the activated sludge and only 3 were biodegraded, resulting in the release of these substances into bodies of water. Textile effluents have been shown to be mutagenic (Claxton et al., 1998; Pinheiro et al., 2004; Sanchez et al., 1988), and mutagenic dyes have been detected in many rivers.

Dyeing is the treatment of fabric with a dye to impart color. Azo dyes are the most important groups. There are main types of dyes used for dyeing different kinds of fibers (Figure 2.1).

Figure 2.1. Dyes for different fibers.

Although the high colouring power of dyes gives rise to aesthetic damage and hence public complaint, several classes of dyes are considered as possible carcinogens or mutagens. It is also noteworthy that some dyes are highly toxic and limiting downstream beneficial uses such as recreation, drinking water and irrigation (Hubbe et al., 2012; Prztstas et al., 2012;). As these dyes are typically very stable in the aquatic environment, the presence of dyes in water bodies, especially those that will be used to produce drinking water, is also a health concern since the population can be exposed to these compounds from the ingestion of both contaminated food and water (Chequer, 2015).

Many dyes do not degrade easily due to their complex structure and textile dye effluent does not decolorize even if the effluent is treated by the municipal wastewater treatment systems (Shaul et al., 1991; Robinson et al., 2002; Forgacs et al., 2004). If textile wastewater, not properly treated is released into the environment, it can introduce metals (Cr and Cd) and organochlorine compounds which can bio-accumulate in fishes in receiving streams.

Before the discharge of textile wastewater into the water bodies, many treatment processes have been developed to treat the wastewater due to the various health risks of the azo dyes to human and aquatic organisms (Figure 2.2.). Of the physical, biological and oxidation methods, ozone treatment which is in the advanced oxidation category, has good decolorization performance (Holkar et. al., 2016). Flocculation is a useful method for the decolorization of wastewater containing disperse dyes although they have low decolorization efficiency for the wastewater having reactive and vat dyes. Contrarily, another physical method, adsorption, has significant decolorization efficiency for wastewater containing a variety of dyes. Oxidation methods are usually used for the degradation of dyes by chemical means due to its easiness of application. The biological methods have some advantages, such as eco-friendly, cost-competitive, less sludge production, giving non-hazardous metabolites and less water consumption compared to physical and oxidation methods (Hayat et al.,

2015). However, textile wastewater parameters after biological treatment are not in compliance with the textile wastewater discharge standards. Therefore, the oxidation process should be applied before biological treatment. On the other hand, if the volume of textile effluent is small, physical and oxidation methods are effective for the removal of dyes in textile wastewater.



Figure 2.2.  Treatment methods applied to textile wastewater.

### 2.1.1. Aminoazobenzene dyes

Aminoazo derivatives are extremely important industrial colorants and widely used in the textile industry (Stead, 1990). As stated above, the azo dyes are carcinogenic, mutagenic or toxic, therefore, they are dangerous to human health as they might release aromatic amines, which in direct and prolonged contact with the human skin or oral cavity. Mathur et al. (2012) reviewed the mutagenicity of azo and non-azo dyes due to extensive recent data on the carcinogenicity and mutagenicity of this group of dyes.

Moreover, it has been proved that variety of 4-aminoazobenzene (AAB), *N*-methyl-4-aminoazobenzene (MAB) and *N,N*-dimethyl-4-aminoazobenzene (DAB) derivatives are mutagenic (Garg et al., 2002).

The main skeleton of aminoazobenzene (AAB) dye and examples of AAB are shown in Figure 2.3 and Figure 2.4, respectively. Freeman et al. (2013) have demonstrated that some small structural modifications of these aminoazobenzene derivatives can reduce or eliminate their mutagenic activity, while maintaining the physical and/or chemical properties that make them useful industrial chemicals. In Figure 2.3, the placement of an azo group between methyl groups produces a colorless compound, while a yellow-orange color is obtained when the azo group is placed between aromatic rings (IARC, 2010). In addition, changing the position of the methoxy group on the phenyl rings dramatically

influences the carcinogenic behavior of the resulting compound. For example, 2-OMe-AAB is noncarcinogenic in rats, whereas 4¢-OMe-AAB is carcinogenic.

Representative chemical structures of some azo dyes and chromophore groups are given in Figure 2.3.



(a)

(b)

(c)

Figure 2.3. Example of aminoazobenzene dye: (a) 4-aminoazobenzene (AAB) and (b) *N*-methyl- 4-aminoazobenzene (MAB) and (c) *N,N*-Dimethyl-4-aminoazobenzene (DAB).

In the past years, several studies were conducted on human exposure to azo dye precursors. Since 70's, intestinal cancer has been more widespread in highly industrialized societies, and thus there may be a relation between the rise in the number of cases of diseases and the use of azo dyes (Wolff and Oehme, 1974; Chung et al., 1978). Later, it has been proven that chronic exposure to aromatic amines would lead to bladder cancer (Bi et al., 1992; Carreon et al., 2010; Sorahan, 2008; You et al., 1990). If an azo compound is ingested orally it can be reduced by anaerobic intestinal microflora and possibly by mammalian azo reductases in the intestinal wall or the liver, to free aromatic amines. Additionally, Zeilmaker et al. (1999) concluded that the dermal exposure to a certain amount of aromatic amines from textiles may lead to cancer risk.

Chung and Cerniglia (1992) reviewed the mutagenicity of azo dyes and found that the MAA, p-phenylenediamine, is a mutagenic moiety of many azo dyes. Lin and Solodar (1988) have also previously pointed out that synthetic azo dyes, which contain a moiety that would be expected to be metabolized to p-phenylenediamine by liver microsomal enzymes, gave positive responses in the Salmonella/microsome assay. It was found that m-diaminobenzene, 2,4-diaminotoluene, 2,4-diaminoethylbenzene are mutagenic to TA98 and TA100 in the presence or absence of metabolic activation (Shahin, 1980). For instance, 2-methoxy-4-aminoazobenzene is an extremely weak mutagen whereas, under similar conditions, 3-methoxy-4-aminoazobenzene is a potent hepatocarcinogen in rats and a strong mutagen in *Escherichia coli* and *Salmonella typhimurium* (Hashimoto et al., 1977; Esancy et al., 1990; Garg et al., 2002; Umbuzeiro et al., 2005a). The potency of these azo dyes is dependent on the position of substituents with respect to both the aromatic rings and the amino nitrogen atom.

Ashby et al. (2006) showed that DAB was not mutagenic, but when pre-incubated with S9, it became mutagenic. They suggested that the mutagenicity of DAB may be attributed to its reductive product *N*,*N*-dimethylphenylenediamine (DMPD). Chung et al. (2006) demonstrated clearly that DMPD was mutagenic with Salmonella tester strain 1538 with S9 microsomal mix.

It should be noted that many of these monocyclic aromatic amines are mutagenic to humans and animals. Therefore, a complete evaluation of the safety of these dyes in the human environment must include an evaluation of their genotoxicity or mutagenicity. It should also be noted that azo dyes which are able to release specific aromatic amines are prohibited in most countries of the world (Chen et al., 2017). The European Community (EC) prevented the production and sale of consumer products dyed with azo derivatives, because azo dyes can be degraded into 22 aromatic amines in the following list in concentrations above 30 ppm in the finished articles or their dyed parts (Table 2.1). EC Directive 2002/61 restricts the use of only about 5% of azo dyes (Cox et al., 2002). Later, the European Commission performed a fast track consultation on a possible restriction of more hazardous substances (CMR 1A & 1B) in textile articles and clothing for consumer use.

Table 2.1. List of amines included in the EC Directive 2002/61

| Aromatic Amine | CAS number |
| --- | --- |
| 4-aminobiphenyl | 92-67-1 |
| Benzidine | 92-87-5 |
| 4-chloro-o-toluidine | 95-69-2 |
| 2-napthylamine | 91-59-8 |
| o-amino-azotoluene | 97-56-3 |
| 5-nitro-o-toluidine | 99-55-8 |
| p-chloroaniline | 106-47-8 |
| 4-methoxy-m-phenylenediamine | 615-05-4 |
| 4,4-methylenedianiline | 101-77-9 |
| 3,3-dichlorobenzidine | 91-94-1 |
| 3,3-dimethoxybenzidine | 119-90-4 |
| 3,3-dimethylbenzidine | 119-93-7 |
| 4,4-methylenedi-o-toluidine | 838-88-0 |
| 6-methoxy-m-toluidine | 120-71-8 |
| 4,4-methyl bis -(2-chloro-aniline) | 101-14-4 |
| 4,4-oxydianiline | 101-80-4 |
| 4,4-thiodianiline | 139-65-1 |
| o-toluidine | 95-53-4 |
| 4-methyl-m-phenylenediamine | 95-80-7 |
| 2,4,5-trimethylaniline | 137-17-7 |
| o-anisidine | 90-04-0 |
| 4-aminoazobenzene | 60-09-3 |

## 2.1.2. Aromatic amine derivatives

Aromatic amines represent one of the most important classes of industrial and environmental chemicals. They have widely used in many industries such as agricultural chemicals, dyes etc. (Woo et al., 2001). Many aromatic amines have also been reported to be carcinogen and mutagen. They are a class of chemicals traditionally recognized as of high concern for human health. They find applications in several chemical industry manufacturing sectors such as oil refining, production of synthetic polymers, adhesives and rubbers, pharmaceuticals, pesticides and explosives (Synyderwine et al., 2002). The structure and numbering system for the parent aromatic amine derivatives used in the present study are given in Figure 2.4.

As early as the late nineteenth century, a doctor related the occurrence of urinary bladder cancer to the occupation of his patients, thus demonstrating concern about the exposure of humans to carcinogenic aromatic amines produced in the dye manufacturing industry. Laboratory investigations subsequently showed that rats and mice exposed to specific azo dye arylamines or their derivatives developed cancer, mainly in the liver (Weisburger, 1997, 2002). In addition, workers in textile dyeing, paper printing, and leather finishing industries, exposed to benzidine based dyes such as Direct Black 38, showed a higher incidence of urinary bladder cancer (Meal et al., 1981).

The large database of mutagenicity results for the aromatic amines has been studied with QSAR approaches (Benigni et al., 2003). Trieff et al. (1989) studied the *Salmonella typhimurium* mutagenicity of 19 aromatic amines tested in the strains TA98 and TA100 with the addition of S9 metabolizing fraction from Aroclor 1254-induced rat liver.



Figure 2.4.  Structure and numbering system for the parent aromatic amine derivatives used in the present study.

Shahin et al. (1983) in a separate study on the mutagenicities of nitro-*p*-phenylenediamine derivatives, discovered that blockage of one amino group in nitro-*p*-phenylenediamine by two hydroxyalkyl groups or blockage of both amino groups, each by one hydroxyalkyl group, eliminated

the mutagenic activity of the compound. In other words, while nitro-*p*-phenylenediamine is a strong mutagen, 3-nitro-4-amino-*N*,*N*-di-hydroxyethyl-aniline, 3- nitro-4-*N*-2-hydroxy-ethyl-amino-*N*-2-hydroxy-ethyl-aniline, 3-nitro-4-*N*-2-aminoethyl-amino-*N*,*N*-dihydroxyethyl-aniline, 3-nitro-4-*N*-methyl-amino-*N*,*N*-dihydroxyethyl-aniline, 3-nitro-4-2-*N*,*N*-diethylamino,ethylamino-*N*,*N*-dihydroxyethyl-aniline, and 3-nitro-4-*N*,*N*-di-ethylamino-propylamino-*N*,*N*-dihydroxyethyl-aniline are not mutagenic. It proves that both free amino groups in the para position and the nitro group are needed for the mutagenic activity of nitro-*p*-phenylenediamine. The mutagenic activity appears to be allied with the position of the amino, nitro, and hydroxy groups in their molecular structures (Shahin, 1985).

Depending on the individual compounds, many aromatic amine metabolites are considered to be non-biodegradable or only very slowly degradable (Saupe, 1999), showing a wide range of toxic effects on aquatic life and higher organisms (Weisburger, 2002; Pinheiro et al., 2004; Khalid et al., 2009).

## 2.2. Ames Test

The Ames test which is an inexpensive and a short-term bacterial reverse mutation assay specifically designed to detect a wide range of chemical substances that can produce genetic damage that leads to gene mutations. Moreover, since the assay is not a live animal model, it fits the 3R principles (Locke, 2006). This test is used worldwide as an initial screen to determine the mutagenic potential of new chemicals and drugs (Mortelmans et al., 2000). The Ames test is an *in vitro* method that commonly uses one of five strains of *Salmonella typhimurium.* In addition*,* the Ames test is by far the most commonly used, long-established *in vitro* test for chemical mutagenicity screening (OECD Test Guideline No.471, 1997).

The *Salmonella* mutagenicity test was specifically designed to detect chemically induced mutagenesis (Mortelmans et al., 2000). The *Salmonella*, or Ames test consists of a range of bacterial strains that together are sensitive to a large array of DNA-damaging agents (Ames, 1984; Zeiger, 1987). TA98 and TA1538 are sensitive to frameshift mutagens, TA100 and TA1535 are used to detect base-pair substitution mutation, TA97 and TA1537 are used for base-pair substitution and some frameshift mutations. TA102 detects mutagens that other strains cannot detect such as formaldehyde. All strains are histidine dependent by virtue of a mutation in the histidine operon (Mortelmans et al., 2000). Histidine is an essential component for growth. Therefore, the bacteria are unable to multiply unless a suitable mutagen causes the proper type of reverse mutation in the histidine gene.

The Ames test with the *Salmonella typhimurium* strains TA98 and TA100 was designed to be compatible with the procedure indicated in the OECD test guideline No. 471, with and without metabolic activation (S9 mix) (OECD Test Guideline No.471, 1997). Test No. 471 is a bacterial reverse mutation test. It uses amino-acid requiring at least five strains of *Salmonella typhimurium* and *Escherichia coli* to detect point mutations by base substitutions or frameshifts. It detects mutations which revert mutations present in the test strains and restore the functional capability of the bacteria to synthesize an essential amino acid (Ames et al., 1975). Harding et al. (2015) could demonstrate the significance of TA98 and TA100 for the detection of aromatic amines mutagenicity. It was noted that the TA98 strain of *Salmonella typhimurium* is sensitive to frameshift mutations (Garg, 2002).

## 2.3. Quantitative Structure-Activity Relationship/ Quantitative Structure-Property Relationship Studies on Dyes

Recently, computational methods have been used to predict the mutagenic potential of azo dyes instead of new experimental studies. Predictive models can be developed to predict the toxicity of azo compounds before their synthesis, based only on their chemical structure. QSAR methods are the most usual alternative for this purpose. They are mathematical models used to predict activities and properties from the physical characteristics (solubility, log Kow, etc.) and/or molecular structure.

In 2006, Registration, Evaluation, Authorisation, and Restriction of Chemicals (REACH) was established by the European Council and the European Parliament. REACH states the need for the evaluation of chemicals that are imported or produced in quantities greater than 1 tonne per annum (tpa) for the assessment of toxic effects by 2018. The European REACH regulation was introduced with the main goal of protecting both human health and the environment. For the assessment of environmental and human hazards, PBT (Persistence, Bioaccumulation, and ecoToxicology) properties should be known as well as CMR (Carcinogenic, Mutagenic and Reproductive toxicity) properties. However, this information is not available for the majority of the existing chemicals.

It is well known that small changes in chemical structure can cause too deep differences in mutagenic activities (Shahin, 1987). The properties of dyes are determined by the structure of the compounds. The potency of azo dyes is strongly dependent on the nature and position of substituents with respect to both the aromatic rings and the amino nitrogen atom (Garg et al., 2002). Substitution of the amino group affects the mutagenicity. For example, 3-nitro-4-*N*-2 hydroxyethylaminoanisole

is mutagenic, 3-nitro-4-*N*-2-hydroxyethyl-aminoaniline weakly mutagenic, but 3- nitro-4-*N*-b-hydroxyethyl-aminophenol is non-mutagenic. In addition, the colour of azobenzene compounds is dependent on the number of –N=N– bonds and the size of the conjugated system (Shahin, 1986). Therefore, studies on the structure-activity relationships may offer the prospect of identifying the critically important parameters in a molecular structure that influence their biological activities like mutagenicity.

Luan et al. (2013) reviewed recent advances and perspectives of QSAR/QSPR studies of dyes. They reported activity/property related research of dyes that have been published in the period from 1995 to 2012. Their emphasis was placed particularly on studies based on QSAR/QSPRs that have contributed to the theoretical design of new, potent and selective dyes. Other studies have concentrated on hazardous properties (such as mutagenicity) as some dyes are toxic, volatile, explosive or radioactive substances. Theoretical methods could provide a means of predicting such properties far more quickly and cheaply than carrying out formal laboratory test protocols. Garg et al. (2002) developed several linear and non-linear QSAR models for the observed mutagenic behaviour of aminoazobenzene derivatives with a variety of molecular descriptors. Three years later, Bhat et al. (2005) developed QSARs that correlated the observed mutagenic activity of 181 aromatic amine derivatives. Common features of the reported QSAR models in these studies are that they do not comply with the OECD principles. According to the Organisation for Economic Co-operation and Development (OECD), a valid QSAR model must have five features. These are (1) a defined endpoint, (2) an unambiguous algorithm, (3) a defined domain of applicability, (4) appropriate measures of goodness of fit, robustness and predictivity and (5) a mechanistic interpretation, if possible (OECD, 2007).

The main steps involved in the development and analysis of a QSAR model can be summarized as: (1) Preparation, analysis, and setup of the input dataset, (2) Model calculation based on selected descriptors, (3) Model exploration, validation, and selection (Roy et al., 2015) (Figure 2.5).

Figure 2.5.  Flowchart of a general QSAR study.

# 3.  MATERIALS AND METHODS

## 3.1. QSAR Workflow

A typical workflow of QSAR used in the present study is given in Figure 3.1. The first stage of the workflow is the selection of a dataset. QSAR modelling was done by following dataset compilation, dataset curation, geometry optimization, and descriptor calculation, dataset splitting, descriptor selection, model selection, and testing the generated QSAR models in terms of some internal and external validation parameters. Then, the best model was tested to achieve predictive ability by using external set compounds.



Figure 3.1. A typical QSAR Workflow.

## 3.2.  Datasets

Two datasets were combined in the present study. The first dataset of mutagenicity data was taken from Garg et al. (2002). This data contains mutagenic activities (rev/nmol) of 43 aminoazobenzene dyes in the TA98 *Salmonella typhimurium* bacterial strain with S9 activation. The second dataset of mutagenic activity data of 181 aromatic amine derivatives was taken from Bhat et

al. (2005). The compounds used in the present study include primary amines, diamines and derivatives of AAB, MAB and DAB. The common feature of this structurally diverse compounds is that they all contain at least one amino group bonded to an aromatic or heteroaromatic ring system (a simple ring or more than one ring forming a conjugated system, fused or nonfused). The 41 compounds of 181 amine derivatives are common in the two datasets. Additionally, the conformers for 6 molecules could not be calculated, because of their high number of conformers. Therefore, 177 compounds were used as a dataset in this study. The mutagenic activity values (rev/nmol) of all chemicals were converted into the logarithmic scale for modelling. As external set chemicals, currently used 33 anionic water-soluble textile dyes (eight anthraquinone and 25 azo dyes) with no experimental mutagenicity data were used. The structure of Acid Yellow 61 (AY61) which is a water-soluble textile dye is given in Figure 3.2.



Figure 3.2.  An example of a water-soluble textile dye: Acid Yellow 61 (AY61).

### 3.3.  Descriptor Calculation

At first, the three-dimensional structure of each molecule was drawn with SPARTAN 10 (Wavefunction Inc., 2010) software. Besides structural drawings, conformational analyses and geometry optimizations were performed at the semi-empirical PM6 level using the same software. An example of the 3D structure of N-hydroxy-2-methoxy-AAB dye drawn in SPARTAN 10 was given in Figure 3.3. For potassium sulfonate ($-SO_3^-K^+$) salt of dye structures were drawn as sulfonic acid form by replacing potassium with hydrogen. The conformer which has the lowest aqueous energy ($E_{aq}$) was saved as .mol2 file. The number of examined conformers depends on number of rotable bonds and flexibility of the molecule. The calculation time for a molecule depends on the number of conformers. One of the conformers of a chemical in the dataset (*N*-hydroxy-2-methoxy-AAB) drawn in SPARTAN software is given in Figure 3.3. After that, these .mol2 files were used for the calculation of molecular descriptors using DRAGON 07 software (Kode, 2017). Type of descriptors calculated by DRAGON 7.0 software is given in Table 2.2.

Table 2.2.  Descriptor blocks and types in DRAGON 7.0 software.

| ID Block | Block Description | Number of Descriptors |
|---|---|---|
| 1 | Constitutional descriptors | 43 |
| 2 | Ring descriptors | 32 |
| 3 | Topological indices | 75 |
| 4 | Walk and path counts | 46 |
| 5 | Connectivity indices | 37 |
| 6 | Information indices | 48 |
| 7 | 2D matrix-based descriptors | 550 |
| 8 | 2D autocorrelations | 213 |
| 9 | Burden eigenvalues | 96 |
| 10 | P_VSA like descriptors | 45 |
| 11 | ETA indices | 23 |
| 12 | Edge adjacency indices | 324 |
| 13 | Geometrical descriptors | 38 |
| 14 | 3D matrix-based descriptors | 90 |
| 15 | 3D autocorrelations | 80 |
| ID Block | Block Description | Number of Descriptors |
| 16 | RDF descriptors | 210 |
| 17 | 3D-MoRSE descriptors | 224 |
| 18 | WHIM descriptors | 114 |
| 19 | GATEWAY descriptors | 273 |
| 29 | Randic molecular profiles | 41 |
| 21 | Functional group counts | 154 |
| 22 | Atom-centered fragments | 115 |
| 23 | Atom-type E-state indices | 170 |
| ID Block | Block Description | Number of Descriptors |
| 24 | CATS 2D | 150 |
| 25 | 2D Atom Pairs | 1596 |
| 26 | 3D Atom Pairs | 36 |
| 27 | Charge descriptors | 15 |
| 28 | Molecular properties | 20 |
| 29 | Drug like indices | 27 |
| 30 | CATS 3D | 300 |

The descriptors generated by the software SPARTAN are: energy of the lowest unoccupied molecular orbital ($E_{LUMO}$), the energy of the highest occupied molecular orbital ($E_{HOMO}$), the energy in the aqueous phase ($E_{aq}$), molecular weight, dipole moments, gas–phase energy ($E$), CPK volume and area (V and A, respectively). Additional descriptors such as ($E_{HOMO}$-$E_{LUMO}$) gap, hardness($\eta$), softness ($\sigma$), chemical potential, electronegativity, and electrophilicity index ($\omega$) were calculated from the energies obtained from SPARTAN 10 according to the equations given below (LoPachin et al., 2007).

$$I(\text{ionization potential}) = -E_{HOMO} \tag{3.1}$$

$$A(\text{electron affinity}) = -E_{LUMO} \tag{3.2}$$

$$X(\text{electronegativity}) = \frac{I+A}{2} \tag{3.3}$$

$$\mu = -X \tag{3.4}$$

$$\omega(\text{electrophilicity index}) = \frac{\mu^2}{2\eta} \tag{3.5}$$

$$\Delta E\text{gap} = E_{HOMO} - E_{LUMO} \tag{3.6}$$

$$\eta(\text{hardness}) = \frac{\Delta Egap}{2} \tag{3.7}$$

$$\sigma(\text{softness}) = \frac{1}{\eta} \tag{3.8}$$

The descriptor values derived from SPARTAN 10 were saved as .txt file to be imported into DRAGON 7.0. The software DRAGON 7.0 can calculate a total 5270 molecular descriptors. The output file of DRAGON software together with SPARTAN 10 derived descriptors were uploaded to QSARINS 2.2.3 (Gramatica et al., 2018) software to develop multiple linear regression (MLR) models for the mutagenic activity (log TA98) of aminoazo dyes (rev/nmol) in *Salmonella typhimurium* (TA98+S9).

SPARTAN is a useful tool to calculate the most stable conformation of compounds and to make geometry optimization. However, it might not be applicative for all chemicals. In this study, we came across some problems for some of the compounds during the calculation of conformer distribution. Thus, for some chemicals whose conformer calculation was not possible with the PM6 method, Molecular Mechanics (MMF) method was employed in SPARTAN 10. Then, we could not calculate aqueous phase energy ($E_{aq}$, whose value is used to determine the lowest energy conformation of a molecule) for some chemicals. At this stage, the gas-phase energy of the compounds was used to select the lowest energy conformer of the molecule. For some compounds, logP value could not be

calculated via SPARTAN, because of the MMF method. Therefore, the calculated logP values from SPARTAN were not included in the descriptor pool.

The output of DRAGON software was saved as a .txt file which was then uploaded into QSARINS (v.2.2.3) software package together with the dependent variable (mutagenic activity) values. Descriptors calculated from SPARTAN 10 were added to this text file as well. Structure drawings, geometry optimization of external set chemicals have been carried out in a previous study (Tugcu et al., 2012). Descriptor values of external set chemicals/structures were calculated similarly.



Figure 3.3.  The conformer of *N*-hydroxy-2-methoxy-AAB drawn in SPARTAN 10 software.

### 3.4.  Training/test set division, descriptor selection and model development

The appearance of some part of the imported data from DRAGON to QSARINS was given in Figure 3.4. Some of the DRAGON descriptors are not calculated for some chemicals, for those descriptor values "na" appears on the relevant column instead of a descriptor value. Therefore, before importing the data from DRAGON file to the software QSARINS, unacceptable descriptors were eliminated as highlighted by the software QSARINS (Figure 3.4).

Figure 3.4. Deletion of column/s with unacceptable data on QSARINs.

QSAR models were developed using the MLR method. MLR is a commonly used method in QSAR studies due to its simplicity, transparency, reproducibility, and easy interpretability. In order to obtain a validated and, therefore, predictive QSAR model, an available dataset should be divided into the training and test sets. Ideally, the division into the training and test sets must satisfy the following three conditions (Golbraikh et al., 2003): (i) All representative compound-points of the test set in the multidimensional descriptor space must be close to those of the training set. (ii) All representative points of the training set must be close to those of the test set. (iii) The representative points of the training set must be distributed within the whole area occupied by the entire dataset. There are many types of division in the literature, such as cluster analysis, periodical division and Kohonen networks (Papa et al., 2005). However, the division of the dataset can also be done randomly. In the present study, we considered three methods for the division of a dataset into training and test sets, (i) ordered by response, (ii) ordered based on the molecular structure, and (iii) random division. We divided the dataset into a training set and a test set comprising 80% and 20% of the whole dataset, respectively.

Response values were ordered in ascending order in both ordered by response and random divisons. These splittings should guarantee that the training set covers the entire range of the experimental responses leaving the minimum and maximum values of response variables (in the

present study, mutagenic activity) in the training set. Divison in the case of ordered by structure, can be done regarding the molecular features of the chemicals. Test set chemicals were not considered in the calculation of the model. It was used only for the external validation to check the efficiency of the prediction of the selected model.

The selection of relevant descriptors, which relate the response variable (mutagenic activity/property/toxicity) to molecular structure, is an important step in the construction of a QSAR model. The most feasible descriptors were selected with "All Subset" and "Hold model add one more variable" tools which are implemented in QSARINS software.

The quality, robustness and predictive power of the model were estimated by different statistical parameters, such as the determination coefficient, either adjusted or not adjusted for the number of variables ($R^2_{adj}$ or $R^2$), the standard deviation (s), the Fisher's statistic ($F$), and the root-mean-square error of training set ($RMSE_{Tr}$). The use of $RMSE_{Tr}$ shows the error between the mean of the experimental values and predicted activities. However, excellent values of $R^2$ and $RMSE_{Tr}$ are not sufficient indicators of a QSAR model validity. Thus, alternative parameters must be provided to indicate the predictive ability of models.

## 3.5. Model Validation

A common method for internally validating a QSAR model is cross-validation (CV, $Q^2$). CV process repeats the regression many times on subsets of data. Usually, each molecule is removed from the dataset once (only), in turn, and the $R$ is computed using the predicted values of the missing molecule. Sometimes more than one molecule (leave many out, LMO) is removed at a time. A cross-validated $R^2$ ($R^2_{CV}$ or $Q^2$) is usually smaller than the overall $R^2$ for a QSAR equation. It is used as a diagnostic tool to evaluate the predictive power of an equation. Cross-validation is used to measure a model's predictive ability and test the over-fitting. Over-fitting refers to the phenomenon in which a predictive model may well describe the relationship between predictors and response, but may subsequently fail to provide valid predictions for new compounds. Over-fitting of a linear QSAR model is usually suspected when the $R^2$ value from the original model is significantly larger (25%) than the $Q^2$ value (Difference between $R^2$ and $Q^2$ should not be more than 0.3) (Leach, 2001). Thus, $Q^2$ is considered a measure of goodness of prediction and not fit in the case of $R^2$. The process of CV begins with the removal of one or a group of compounds, which becomes a temporary test set, from the training set. In the leave-one-out (LOO) method of CV, the process of removing a molecule, and creating and validating the model against the individual molecules is performed for the entire training

set. Once complete, the mean is taken of all the $Q^2$ values and reported. According to the rule-of-thumb, the ratio of the number of compounds in the training set over the number of variables (descriptors) should have a value of at least 5 which is called as Topliss ratio suggesting that at least 5 training compounds should be represented with one descriptor. For instance, a training set that has 25 training set compounds should not have more than 5 descriptors (Topliss and Edwards, 1979). Inter-correlation among descriptors was tested via the QUIK (Q Under Influence of K) rule (Todeschini et al., 2004). QUIK rule was set to 0.05 to minimize inter-correlation between descriptors.

The predictive power of the equation is poor when the observations are not sufficiently independent of each other. One way to test for this is by randomization of the dependent variable (response/activity). This procedure ensures that the model is not due to a chance. Activity values are shuffled randomly and the entire modelling procedure is repeated. This process is repeated many times. The significantly low coefficients of determination of new models indicate that the proposed model is not obtained by chance correlation. Therefore, the reliability of the models was tested by response randomization (*Y*-scrambling) procedure.

A developed QSAR model is accepted if it can satisfy the following criterion (The following values are the minimum recommended values for significant QSAR model (Veerasamy et al., 2011).

- If coefficient of determination $R^2 > 0.6$
- If the standard deviation *s* is not much larger than the standard deviation of the activity data.
- If its *F* value indicates that the overall significance level is better than 95%.
- If its confidence interval of all individual regression coefficients proves that they are justified at the 95% significance level.

Golbraikh and Tropsha (2003) identified some criteria for external prediction. If the following criteria are valid, the models are acceptable:

(a)  If cross-validated $R^2_{CV}$ $(Q^2) > 0.5$

(b)  If $R^2$ for external/test set, $R^2_{Test} > 0.6$

(c)  $(R^2 - R^2_0)/R^2 < 0.1$ and $0.85 < k < 1.15$

(d)  $(R^2 - R'^2_0)/R^2 < 0.1$ and $0.85 < k' < 1.15$ (for test set).

(e)  $r^2_m$ (overall) $> 0.5$ (or at least near 0.5); $\Delta r^2_m < 0.2$.

where $R^2_0$ (predicted vs. observed) and $R'^2_0$ (observed vs. predicted) are the determination coefficients without intercept, $k$ and $k'$ are the slopes. The parameter $r^2_m$ (overall) penalizes a model for large differences between observed and predicted values of the compounds of the whole set (considering both training and test sets). $\Delta r^2_m$ estimated the closeness between the values of the predicted and the corresponding observed activity data.

In addition, the response data (mutagenic activity) should cover a range of at least two or even more logarithmic units. They should be well distributed over the whole distance.

An equation has to be rejected:
- If the above mentioned statistical measures are not satisfied
- If the number of the variables in the regression equation is unreasonably large.
- If the standard deviation is smaller than the error in the activity data.

Additional external validation criteria to be used in modelling are: $Q^2_{F1}$, $Q^2_{F2}$ (Shüürman et al., 2008) $Q^2_{F3}$ (Consonni et al., 2009, 2010), Concordance Correlation Coefficient (CCC) (Lin 1989, 1992), for training and test sets $CCC_{Tr}$ and $CCC_{Test}$ (Chirico and Gramatica 2011, 2012), respectively. $Q^2_{F1}$ shows the degree of correlation between the experimental and predicted activity of the dataset (Shi et al., 2001). The $Q^2_{F2}$ parameter was described by Schüürmann et al. (2008). The main difference between $Q^2_{F1}$ and $Q^2_{F2}$ is that the mean experimental activity is replaced in $Q^2_{F2}$ with the mean predicted activity. $Q^2_{F3}$ measures the model's predictivity it is sensitive to training set selection and it criticizes the dataset when they are very homogeneous. Concordance Correlation Coefficient ($CCC$) measures both the distance of observations to the fitting line and the distance which the regression line deviates from slope 1 passing through the origin. Thus, $CCC$ value is often smaller than its ideal value of 1 (Chirico and Gramatica, 2011). The external validation ensures the predictability and applicability of the developed QSAR model for the prediction of untested molecules. In the present study, current textile dyes were used as an external set.

Multi-Criteria Decision Making (MCDM) implemented in the software QSARINS was used for the ranking of generated models. It consists of the summary of performances of a certain number of criteria associated with both the internal and external validation. Each validation criteria value ranges from 0 to 1 and the geometric average of all the values obtained from the desirability functions creates the MCDM value (QSARINS 2.2.3). The MCDM of fitting (maximizing $R^2$, $R^2_{adj}$, and $CCC_{Tr}$ while minimizing $R^2 - R^2_{adj}$) cross-validation (maximizing $Q^2_{LOO}$, $Q^2_{LMO}$ and $CCC_{CV}$, while minimizing $R^2_{Y\text{-SCR}}$) and external validation (maximizing $Q^2_{F1}$, $Q^2_{F2}$, $Q^2_{F3}$, and $CCC_{Test}$) parameters are

automatically calculated using all the corresponding criteria. The model with the best MCDM compromise among the selected validation criteria is sorted as the best. Additionally, for the model selection, the principle of parsimony was taken into account; the model with the highest statistical significance (both model and test set significance), but having as few parameters as possible was selected.

Besides the above-mentioned criteria, Mean Absolute Error (*MAE*)-based criteria proposed by Roy et al., (2016) was applied to test further the external validation of all models selected by MCDM. *MAE* provides equal weight to all the errors (Roy et al., 2016). Xternal Validation Plus tool developed by Roy et al. (2015) was used to check the presence of systematic errors in the model. It further computes all the required external validation parameters, while judges the performance of actual prediction quality of a QSAR model based on the mean absolute error. The mean absolute error (*MAE*) can be calculated with the following equation:

$$MAE = \frac{\sum_{i=1}^{n_{ext}} |y_i - \hat{y}_i|}{n_{ext}} \qquad (3.9)$$

I. Good predictions:

From a general notation, an error of 10% of the training set range should be acceptable while an error value more than 20% of the training set should be a very high error. Thus, the criteria for good predictions should be the following:

*MAE* ≤ 0.1 x training set range and *MAE* + 3σ ≤ 0.2 x training set range.

Where, the σ value refers to the standard deviation of the absolute error values for the test set data. Considering a normal distribution pattern, mean ± 3σ covers 99.7% of the data points.

II. Bad predictions:

A value of *MAE* more than 15% of the training set range should be high while an error of more than 25% of the training set is considered very high. Hence, the predictions could be considered very high. Hence, the predictions could be considered when:

*MAE* > 0.15 x training set range or *MAE* + 3σ > 0.25 x training set range.

The predictions which do not fall under either of the above two conditions may be considered as of moderate quality. The mentioned criteria should be used in cases where there are more than 10 data points in the test set (Roy et al., 2016). If anyone or more conditions stated in Figure 3.5. were met, a systematic error occurs. If the systematic error present in the output file, the model should be discarded.

(i) NPE/NNE > 5* or NNE/NPE >5*

Here, **NPE:** Number of positive errors AND **NNE:** Number of negative errors

(ii) ABS(MPE/MNE) > 2* or ABS(MNE/MPE) > 2*

Here, **MPE:** Mean of positive error; **MNE:** Mean of negative errors; **ABS:** Absolute value

(iii) AAE - ABS(AE) < 0.5*×AAE

Here, **AAE:** Average of absolute prediction errors; **AE:** Average of predictions errors; **ABS:** Absolute value

(iv) $R^2$ ($i^{th}$ vs $(i-1)^{th}$ residuals) > 0.5* for residuals sorted on experimental response values (Yobs)

(v) $R^2$(Yobs and residuals) > 0.5*

Figure 3.5. Error types of Xternal validation plus tool.

"Prediction Reliability Indicator" was developed by Roy et al., (2018) to indicate or categorize the quality of predictions for the test set (known experimental response) or external (unknown experimental response) sets into three groups: good (with composite score 3), moderate (with composite score 2) and bad (with composite score 1). There is no clear discussion provided for the output of this program. For instance, it tabulates the prediction of the response variable of a chemical from the model as good, but it puts a mark for that chemical as out of the applicability domain (AD) of model. Therefore, we did not apply PRI to our selected models.

## 3.6. Applicability Domain of QSAR Model

The derivation of QSAR models is based primarily on training sets which are structurally limited and thus, their applicability to the query chemicals is limited (Dimitrov et al., 2005). Thus, their applicability towards reliable predictions is restricted in chemical space to some specific categories. The applicability domain (AD) is a theoretical region of the chemical space, defined by the model descriptors and modelled response and, thus, by the nature of the training set molecules. It is important to note that the AD of a model cannot be verified by studying only a few chemicals (even less than five), as in such cases, it could happen that extrapolated predictions are good, but probably only by chance, so it is impossible to rely on the possibility of obtaining general conclusions. AD is one of

the five OECD criteria (OECD, 2007). It estimates the similarity of an individual compound's to the rest of the dataset.

The AD of a QSAR should be described in terms of the most relevant parameters i.e. usually those that are the range of descriptors and response.

In the present study, the AD of the generated MLR model was analyzed by Williams plot. To check whether a new chemical lies within the applicability domain, the leverage approach was used. A compound was considered outside the AD when the leverage value is higher than the warning limit of $3p/n$, where p is the number of model variables plus 1 and $n$ is the number of objects used to develop the model. This indicates structural ($X$) outliers. Compounds with a residual higher than $\pm 3\sigma$ in the training set were considered as response ($Y$) outliers. Additionally, Insubria Graph was used to predict the unknown endpoints of chemicals from the model equation regarding the applicability domain of the generated QSAR model. In the present study, the predicted ability of all models was tested using an external dataset which comprises currently used 33 anionic water-soluble textile dyes (eight anthraquinone and 25 azo dyes) with no measured mutagenic activity. The chemical structures of these dyes are given in Appendix A.

# 4. RESULTS AND DISCUSSION

## 4.1. Model Development

The compiled experimental mutagenicity data from literature were given in Table 4.1 together with the molecular weight and logP values obtained from Garg et al., 2002 and Bhat et.al., 2005). The normality of the mutagenicity data was evaluated using the Kolmogorov-Smirnov test in SPSS 25 (IBM). The distribution of data was found to be normal ($p > 0.05$). Fortunately, the mutagenic activity of compounds ranges over nearly 7 orders of magnitude, providing a broad range of values for QSAR development.

The dataset contained 177 chemicals was divided into two groups as test and training sets. The training set was composed of 80% of the whole data with 142 compounds and the test set was composed of 20% of the whole data with 35 compounds. Test set compounds in each division and the ratio of the number of test set to the training set compounds in each division are given in Table 4.2.

Numerous training/test divisions were tried and many different 3 to 6-descriptor QSAR models were generated using "AllSubset" and "Hold model add one more variable" tools of QSARINS (v.2.2.3) software. The generated models ranked by MCDM as well as the internal and external validation criteria. Fit, internal and external validation parameters of the best 3 to 6 descriptor models are given in Table 4.3 and 4.4, respectively. The generated models fulfilled all the fit, internal and external validation criteria, therefore they are subjected to further criteria to determine the best model. As such, all models were tested for their external prediction capacity regarding the highest number of compounds from external set chemicals which fell within the applicability domain. Models from each division were selected regarding the highest $R^2$ and $Q^2$ values, the minimum number of structural outliers and higher predictive performance of the external set chemicals and listed in Table 4.3.

The internal predictive power of each model is judged by the parameters like $R^2$ (the squared of determination coefficient) and $Q^2_{LOO}$ (the leave-one-out cross-validation). For the 3-6-descriptor models, $R^2$ and $Q^2_{LOO}$ values were very similar to each other which reveal the stability of all models ($R^2$ range is between 0.54 and 0.61, and $Q^2_{LOO}$ range is between 0.51 and 0.57). The models' reliability and robustness were checked by the *Y*-scrambling procedure. The mutagenicity values were shuffled

randomly and new regression models were developed. The generated models' coefficients of determination were calculated. The procedure was repeated 2000 times. The average $R^2$ of shuffled models were significantly low (results were not shown here), revealing that there is no chance correlation in all models.

The high $R^2_{Test}$ values (0.74 to 0.77) and low $RMSE_{Test}$ relative to the $RMSE_{Tr}$ proved to have good predictive power for all models. Additionally, all models passed the external validation criteria mentioned in section 3.5 and threshold values reported in the literature. All these results indicate that the developed models are robust, validated and predictive. For all chemicals in the training and the test sets, the standardized residuals are smaller than three standard deviation units (Williams plot). All models have no response outlier. The predicted vs. observed mutagenic activity and Williams plots of all models except model M2 are given in Figure B1 and Figure B2, respectively (Appendix B).

The structural coverage of all models for external set chemicals ranged from 58% to 70% (Table 4.5). Of the five models, all modelwere classified as GOOD regarding the $MAE$-based criteria.

Table 4.1. The name, LogP, molecular weight and mutagenic activity values of chemicals in the dataset.



Main Skeleton of Compond ID 1-42          Main Skeleton of Compound ID 162-177

| ID | Compound Name | log P | Molecular Weight (g/mol) | Mutagenic activity in TA98+S9 (rev/nmol) | CAS Number |
|---|---|---|---|---|---|
| 1 | 4'-NEt2-3-OMe-AAB | 5.16 | 298.39 | -2.15 | **- |
| 2 | 2-OMe-AAB | 3.87 | 227.27 | -2.00 | 80830-39-3 |
| 3 | 4'-OH-AAB | 2.55 | 213.24 | -1.28 | 6530-27-4 |
| 4 | 3'-Me-4'-OH-AAB | 3.01 | 227.27 | -2.23 | - |
| 5 | 4'-OH-2',3-diMe-AAB(4'-OH-OAT) | 3.47 | 241.29 | -0.95 | - |
| 6 | AAB | 3.13 | 197.24 | -0.69 | 60-09-3 |
| 7 | 3'-Me-AAB | 3.59 | 211.27 | -0.62 | - |
| 8 | 3-OMe-4'-N(CH2CH2OH)2-AAB | 2.58 | 330.39 | -0.41 | - |
| 9 | 3'CH2OH-AAB | 1.94 | 227.27 | -0.22 | - |
| 10 | 3-OH-AAB | 2.97 | 213.24 | -0.16 | 14528-97-3 |
| 11 | 3-OCH2CH2OH-AAB | 2.51 | 257.29 | 0.13 | - |
| 12 | 2'-CH2OH-3-Me-AAB | 2.40 | 241.29 | 0.30 | - |
| 13 | 4-OMe-AAB | 2.95 | 227.27 | 0.36 | - |
| 14 | 2',3-diMe-AAB | 4.05 | 225.30 | 0.43 | 97-56-3 |
| 15 | 3-OBu-AAB | 5.08 | 269.35 | 0.70 | - |
| 16 | 3-OEt-AAB | 4.02 | 241.29 | 1.14 | 126335-27-1 |
| 17 | 3-OPr-AAB | 4.55 | 255.32 | 1.28 | - |
| 18 | 3-OMe-AAB | 3.48 | 227.27 | 1.89 | - |
| 19 | 3'-Me-4'-OH-MAB | 3.67 | 241.29 | -1.15 | - |
| 20 | 3'-COOH-MAB | 3.52 | 255.28 | -0.91 | - |
| 21 | 4'-OH-MAB | 3.21 | 227.27 | -0.85 | - |
| 22 | MAB | 3.79 | 211.27 | -0.74 | 621-90-9 |
| 23 | 4'-Me-MAB | 4.25 | 225.30 | -0.55 | 28149-22-6 |
| 24 | 3'-Me-MAB | 4.25 | 225.30 | -0.35 | 2058-62-0 |
| 25 | 3'-CH2OH-MAB | 2.60 | 241.29 | -0.30 | - |
| 26 | 3'-Me-4'-OH-DAB | 4.31 | 255.32 | -0.96 | - |
| 27 | p-(dimethylamino)azobenzene | 4.43 | 225.30 | -0.85 | 60-11-7 |
| 28 | 3'-COOH-DAB | 4.17 | 269.30 | -0.70 | - |
| 29 | 2-Me-DAB | 4.89 | 239.32 | -0.66 | 54-88-6 |

Table 4.1.  Continued

| ID | Compound Name | log P | Molecular Weight (g/mol) | Mutagenic activity in TA98+S9 (rev/nmol) | CAS Number |
|---|---|---|---|---|---|
| 30 | 3'-Me-DAB | 4.89 | 239.32 | -0.45 | 55-80-1 |
| 31 | 3'-CHO-DAB | 3.91 | 253.31 | -0.42 | - |
| 32 | 3'-CH2-OAC-DAB | 4.14 | 297.36 | 2.71 | - |
| 33 | 3'-CH2OH-DAB | 3.25 | 255.32 | -0.22 | 35282-68-9 |
| 34 | 3'-Me-AAB-N-Ac | 3.73 | 253.31 | -1.06 | - |
| 35 | 3'-Me-4'-OH-AAB-N-Ac | 3.15 | 269.30 | -1.05 | - |
| 36 | N-OH-2-OMe-AAB | 4.08 | 243.27 | -0.96 | - |
| 37 | 3'-Me-MAB-N-Ac | 3.03 | 267.33 | -0.28 | - |
| 38 | N-OH-MAB | 2.74 | 227.27 | -0.19 | 1910-36-7 |
| 39 | N-OH-3'-Me-MAB | 3.20 | 241.29 | 0.00 | - |
| 40 | N-OH-AAB | 2.98 | 213.24 | 0.01 | 6530-27-4 |
| 41 | N-OH-4'-Me-MAB | 3.20 | 227.27 | 0.05 | - |
| 42 | N-OH-3-OMe-AAB | 3.08 | 243.27 | 2.28 | - |
| 43 | 2-aminobiphenyl | 2.68 | 169.23 | *-1.49 | 90-41-5 |
| 44 | 4-aminobiphenyl | 2.77 | 169.23 | -0.14 | 92-67-1 |
| 45 | 2-2'-diaminobiphenyl | 1.39 | 184.24 | -1.52 | 1454-80-4 |
| 46 | 2,4'-diaminobiphenyl | 1.47 | 184.24 | -0.92 | 492-17-1 |
| 47 | 3,3'-diamnobiphenyl | 1.41 | 184.24 | -1.30 | 2050-89-7 |
| 48 | 3,4'-diaminobiphenyl | 1.48 | 184.24 | 0.20 | 32316-90-8 |
| 49 | 2-amino-3'-nitrobiphenyl | 2.21 | 216.24 | -0.89 | 34862-87-8 |
| 50 | 2-amino-4'nitrobiphenyl | 2.30 | 214.22 | -0.62 | 6272-52-2 |
| 51 | 3-amino-2'nitrobiphenyl | 2.24 | 214.22 | -1.30 | 96187-18-7 |
| 52 | 3-amino-3'nitrobiphenyl | 2.22 | 214.22 | -0.55 | 31835-64-0 |
| 53 | 3-amino-4'nitrobiphenyl | 2.31 | 214.22 | 0.69 | 53059-29-3 |
| 54 | 4-amino-2'nitrobiphenyl | 2.32 | 214.22 | -0.92 | 1140-28-9 |
| 55 | 4-amino-3'nitrobiphenyl | 2.30 | 214.22 | 1.02 | 1141-29-3 |
| 56 | 4-amino-4'nitrobiphenyl | 2.39 | 214.22 | 1.04 | 1211-40-1 |
| 57 | 1-aminoanthracene | 3.40 | 193.25 | 1.18 | 610-49-1 |
| 58 | 2-aminoanthracene | 3.40 | 193.25 | 2.62 | 613-13-8 |
| 59 | 9-aminoanthracene | 3.40 | 193.25 | 0.87 | 779-03-3 |
| 60 | 1-aminophenanthrene | 3.40 | 193.25 | 2.38 | 4176-53-8 |
| 61 | 2-aminophenanthrene | 3.40 | 193.25 | 2.46 | 3366-65-2 |
| 62 | 3-aminophenanthrene | 3.40 | 193.25 | 3.77 | 1892-54-2 |
| 63 | 9-aminophenanthrene | 3.40 | 193.25 | 2.98 | 947-73-9 |
| 64 | 1-aminofluorene | 2.88 | 181.24 | 0.43 | 6344-63-4 |
| 65 | 2-aminofluorene | 2.88 | 181.24 | 1.93 | 153-78-6 |
| 66 | 3-aminofluorene | 2.88 | 181.24 | 0.89 | 6344-66-7 |
| 67 | 4-aminofluorene | 2.88 | 181.24 | 1.13 | 7083-63-8 |
| 68 | 2,7-diaminofluorene | 1.60 | 196.25 | 0.48 | 525-64-4 |
| 69 | 2-amino-7-nitrofluorene | 2.61 | 226.24 | 3.00 | 1214-32-0 |

Table 4.1.  Continued

| ID | Compound Name | log P | Molecular Weight (g/mol) | Mutagenic activity in TA98+S9 (rev/nmol) | CAS Number |
|---|---|---|---|---|---|
| 70 | 2-bromo-7-aminofluorene | 3.65 | 260.13 | 2.62 | 6638-60-4 |
| 71 | 2-hydroxy-7-aminofluorene | 2.15 | 197.24 | 0.41 | 1953-38-4 |
| 72 | 2-amino-7-acetamidofluorene | 1.75 | 238.29 | 1.18 | 6957-50-2 |
| 73 | 1-aminopyrene | 4.31 | 217.27 | 1.43 | 1606-67-3 |
| 74 | 2-aminopyrene | 3.81 | 217.27 | 3.50 | 1732-23-6 |
| 75 | 4-aminopyrene | 3.89 | 217.27 | 3.16 | 17075-03-5 |
| 76 | 1-aminofluoranthene | 3.89 | 217.27 | 3.35 | 13177-25-8 |
| 77 | 2-aminofluoranthene | 3.89 | 217.27 | 3.23 | 13177-26-9 |
| 78 | 3-aminofluoranthene | 3.89 | 217.27 | 3.31 | 2693-46-1 |
| 79 | 7-aminofluoranthene | 3.89 | 217.27 | 2.88 | 5869-25-0 |
| 80 | 8-aminofluoranthene | 3.89 | 217.27 | 3.80 | 2642-98-0 |
| 81 | 6-aminochrysene | 4.63 | 243.31 | 1.83 | 580-17-6 |
| 82 | 3-aminoquinoline | 1.49 | 144.18 | -3.14 | 611-34-7 |
| 83 | 5-aminoquinoline | 0.92 | 144.18 | -2.00 | 580-15-4 |
| 84 | 6-aminoquinoline | 1.26 | 144.18 | -2.67 | 578-66-5 |
| 85 | 8-aminoquinoline | 1.88 | 144.18 | -1.14 | 18992-86-4 |
| 86 | 1-aminocarbazole | 2.43 | 182.23 | -1.04 | 4539-51-9 |
| 87 | 2-aminocarbazole | 2.43 | 182.23 | 0.60 | 6377-12-4 |
| 88 | 3-aminocarbazole | 2.44 | 182.23 | -0.48 | 18992-64-8 |
| 89 | 4-aminocarbazole | 2.44 | 182.23 | -1.42 | 2876-22-4 |
| 90 | 1-aminophenazine | 3.19 | 195.23 | -0.01 | 2876-23-5 |
| 91 | 2-aminophenazine | 2.13 | 195.23 | 0.55 | 16582-03-9 |
| 92 | 1,6-diaminophenazine | 3.53 | 210.24 | 0.20 | 28124-29-0 |
| 93 | 1,7-diaminophenazine | 3.53 | 210.24 | 0.75 | 102877-14-5 |
| 94 | 1,9-diaminophenazine | 3.53 | 210.24 | 0.04 | 120209-97-4 |
| 95 | 2,7-diaminophenazine | 1.64 | 210.24 | 3.97 | 7704-40-7 |
| 96 | 2,8-diaminophenazine | 1.54 | 210.24 | 1.12 | 134-32-7 |
| 97 | 1-naphthylamine | 2.17 | 143.19 | -0.60 | 91-59-8 |
| 98 | 2-naphthylamine | 2.17 | 143.19 | -0.67 | 776-34-1 |
| 99 | 1-amino-4-nitronaphthalene | 2.62 | 188.19 | -1.77 | 606-57-5 |
| 100 | 2-amino-1-nitronaphthalene | 3.06 | 188.19 | -1.17 | 62-53-3 |
| 101 | aniline | 0.93 | 93.13 | -1.51 | 106-40-1 |
| 102 | *p*-bromoaniline | 2.05 | 172.03 | -2.70 | 95-51-2 |
| 103 | *o*-chloroaniline | 1.91 | 127.57 | -3.00 | 106-47-8 |
| 104 | *p*-chloroaniline | 1.76 | 127.57 | -2.52 | 371-40-4 |
| 105 | *p*-fluoroaniline | 1.15 | 111.12 | -3.32 | 156-43-4 |
| 106 | *p*-ethoxyaniline | 1.28 | 137.18 | -2.30 | 139-59-3 |
| 107 | 4-phenoxyaniline | 2.77 | 185.23 | 0.38 | 123-30-8 |
| 108 | 4-hydroxyaniline | -0.30 | 109.13 | -1.60 | 6373-50-8 |
| 109 | 4-cyclohexylaniline | 3.46 | 175.28 | -1.24 | 97-02-9 |

Table 4.1.  Continued

| ID | Compound Name | log P | Molecular Weight (g/mol) | Mutagenic activity in TA98+S9 (rev/nmol) | CAS Number |
|---|---|---|---|---|---|
| 110 | 2,4-dinitroaniline | 2.22 | 183.12 | -2.00 | 95-68-1 |
| 111 | 2,4-dimethylaniline | 1.86 | 121.18 | -2.22 | 95-78-3 |
| 112 | 2,5-dimethylaniline | 1.86 | 121.18 | -2.40 | 367-25-9 |
| 113 | 2,4-difluoroaniline | 1.49 | 129.11 | -2.70 | 120-71-8 |
| 114 | 2-methoxy-5-methylaniline | 1.55 | 137.18 | -2.05 | 16452-01-0 |
| 115 | 3-methoxy-4-methylaniline | 1.29 | 137.18 | -1.96 | 102-50-1 |
| 116 | 4-methoxy-2-methylaniline | 1.20 | 137.18 | -3.00 | 89-63-4 |
| 117 | 4-chloro-2-nitroaniline | 2.74 | 172.57 | -2.22 | 137-17-7 |
| 118 | 2,4,5-trimethylanilne | 2.32 | 135.21 | -1.32 | 1817-73-8 |
| 119 | 2-bromo-4,6-dinitroaniline | 3.09 | 262.02 | -0.54 | 621-95-4 |
| 120 | 4,4'-ethylenebis(aniline) | 2.27 | 212.30 | -2.15 | 101-77-9 |
| 121 | 4,4'-methylenedianiline | 1.64 | 198.27 | -1.60 | 578-54-1 |
| 122 | 4,4'-methylenebis (O-ethylaniline) | 2.32 | 121.18 | -0.99 | 348-54-9 |
| 123 | 4,4'-methylenebis(o-fluoroaniline) | 2.50 | 111.12 | 0.23 | 95-54-5 |
| 124 | 4,4'-methylenebis(o-isopropylaniline) | 4.31 | 135.21 | -1.77 | 108-45-2 |
| 125 | 1,2-Phenylenediamine | 0.05 | 108.14 | -0.75 | 106-50-3 |
| 126 | 1,3-benzenediamine | -0.82 | 108.14 | -0.46 | 95-70-5 |
| 127 | 1,4-benzenediamine | -0.90 | 108.14 | -0.89 | 5307-14-2 |
| 128 | 2-methyl-1,4-phenylenediamine | -0.45 | 122.17 | -1.52 | 5307-02-8 |
| 129 | 2-nitro-1,4-phenylenediamine | 0.75 | 153.14 | -0.05 | 90-41-5 |
| 130 | 2-methoxy-1,4-phenylenediamine | -0.74 | 138.17 | 0.32 | 92-67-1 |
| 131 | 2-ethoxy-1,4-phenylenediamine | -0.21 | 152.20 | -0.02 | - |
| 132 | 2-propoxy-1,4-phenylenediamine | 0.29 | 166.22 | -0.21 | - |
| 133 | 2-butyloxy-1,4-phenylenediamine | 0.82 | 180.25 | -3.00 | - |
| 134 | 4-chloro-1,2-phenylenediamine | 1.15 | 142.59 | -0.49 | 95-83-0 |
| 135 | 4-chloro-1,3-phenylenediamine | 0.38 | 142.59 | -0.77 | 5131-60-2 |
| 136 | 4-nitro-1,2-phenylenediamine | 1.21 | 153.14 | 0.35 | 99-56-9 |
| 137 | 4-nitro-1,3-phenylenediamine | 0.80 | 153.14 | -2.40 | 5131-58-8 |
| 138 | *N*-acetyl-1,4-phenylenediamine | 0.08 | 150.18 | -1.80 | 122-80-5 |
| 139 | 1,4-*N*Ac-phenylenediamine | -0.45 | 192.22 | -1.43 | 140-50-1 |
| 140 | 2,6-dichloro-1,4-phenylenediamine | 1.68 | 177.03 | -0.69 | 609-20-1 |
| 141 | *N,N*-diethyl-1,4-phenylenediamine | 0.32 | 164.25 | -2.15 | 93-05-0 |
| 142 | *N,N*-dimethyl-1,4-phenylenediamine | -0.58 | 136.20 | -0.87 | 99-98-9 |
| 143 | *N*-methyl-1,4-phenylenediamine | -0.39 | 122.17 | -0.38 | 623-09-6 |
| 144 | 2,4-diaminoethylbenzene | 0.17 | 136.20 | -0.87 | 1195-06-8 |
| 145 | 2,4-diaminotoluene | -0.36 | 122.17 | -1.29 | 95-80-7 |
| 146 | 3,4-diaminotoluene | 0.51 | 122.17 | -1.42 | 496-72-0 |
| 147 | 3-amino-α,α,α-trifluorotoluene | 2.30 | 161.13 | -0.80 | 98-16-8 |

Table 4.1.  Continued

| ID | Compound Name | log P | Molecular Weight (g/mol) | Mutagenic activity in TA98+S9 (rev/nmol) | CAS Number |
|---|---|---|---|---|---|
| 148 | 2,4-diamino-isopropylbenzene | 0.52 | 150.23 | -3.00 | 14235-45-1 |
| 149 | 2,4-diamino-*n*-butylbenzene | 1.24 | 164.25 | -2.70 | 63921-07-3 |
| 150 | 2-amino-4-chlorophenol | 1.67 | 143.57 | -3.00 | 95-85-2 |
| 151 | 2-amino-4-methylphenol | 0.90 | 123.16 | -2.10 | 95-84-1 |
| 152 | 2-amino-5-nitrophenol | 1.61 | 154.13 | -2.52 | 121-88-0 |
| 153 | 3-amino-6-methyl phenol | 0.80 | 123.16 | -1.40 | 2835-95-2 |
| 154 | benzidine | 1.56 | 184.24 | -0.39 | 92-87-5 |
| 155 | 3,3'-diaminobenzidine | -0.21 | 214.27 | -0.04 | 91-95-2 |
| 156 | 3,3'-dichlorobenzidine | 3.45 | 253.13 | 0.81 | 91-94-1 |
| 157 | 4,4'-diamino-3,3'-dimethoxy-1,1'-biphenyl | 1.58 | 244.29 | 0.15 | 119-90-4 |
| 158 | 3,3'-dimethyl-4,4'-bianiline | 2.48 | 212.30 | 0.01 | 119-93-7 |
| 159 | 4-aminophenyl disulfide | 1.65 | 248.37 | -1.03 | 722-27-0 |
| 160 | (4-aminophenyl)ether' | 0.51 | 200.24 | -1.14 | 287476-22-6 |
| 161 | 4-aminophenyl sulfide | 3.07 | 216.31 | 0.31 | 139-65-1 |
| 162 | 3-OSO$_3$K-AAB | -1.00 | 293.30 | -1.43 | - |
| 163 | 4'-OSO$_3$K-AAB | -1.29 | 293.30 | -1.01 | - |
| 164 | 3-methoxy-4'-N,N-diethyl-AAB | 5.16 | 313.41 | -2.15 | - |
| 165 | 3-OSO$_3$K-MAB | -1.05 | 307.33 | -1.47 | - |
| 166 | 4'-OSO$_3$K-MAB | -1.33 | 307.33 | -0.38 | - |
| 167 | 3'-acetoxymethyl-DAB | 4.14 | 297.36 | -0.29 | - |
| 168 | R1=OCH$_2$CH$_2$OH, R2=NEt2 | 0.12 | 562.65 | -0.83 | - |
| 169 | R1=OBu, R2=NEt2 | 2.69 | 574.71 | -0.64 | - |
| 170 | R1=OPr, R2=NEt2 | 2.16 | 560.68 | -0.59 | - |
| 171 | R1=OMe, R2=NEt2 | 1.09 | 532.63 | -0.47 | - |
| 172 | R1=OPr, R2=H | 0.39 | 489.56 | -0.40 | - |
| 173 | R1=OEt, R2=NEt2 | 1.63 | 546.65 | -0.24 | - |
| 174 | R1=OBu, R2=H | 0.92 | 503.58 | -0.18 | - |
| 175 | R1=OCH$_2$CH$_2$OH, R2= H | -1.65 | 475.53 | 0.13 | - |
| 176 | R1=OMe, R2=H | -0.67 | 461.50 | 0.47 | - |
| 177 | R1=OEt, R2=H | -0.14 | 475.53 | 0.64 | - |

*The mutagenic activity values of chemicals with ID values between 1-43 were taken  from Garg et al., 2002 ; the others were taken from Bhat et al., 2005.

**data could not be found

Table 4.2.  Test set compounds for three divisions and five models used in mutagenicity modelling.

| Model | $n_{Test}/n_{Tr}$ | Test Set Compounds* |
|---|---|---|
| **Division 1** | | |
| M1 | 35/142 | 3,4,7,8,12,15,29,38,39,41,50,57,59,61,65,66,85,86,93,98, |
| M2 | 35/142 | 101,108,114,116,126,134,138,139,148,150,154,155,156,160,174 |
| **Division 2** | | |
| M3 | 34/143 | 1,3,13,21,27,36,37,39,40,41,45,48,54,56,59,70,71,72,79,87,90,92, 100,109,114,115,118,119,124,127,141,146,153,159 |
| **Division 3** | | |
| M4 | 34/143 | 1,3,13,21,27,36,37,39,40,41,45,48,54,56,59,70,71,72,79,87,90,92, 100,109,114,115,118,119,124,128,141,146,153,159 |
| M5 | 34/143 | |

**\*** Compound numbers refer to the ID numbers given in Table 4.1

Table 4.3. The fit and internal parameters of the developed QSAR models for the mutagenic activity of textile dyes

| Model No | Number of Variables | Variables | Fitting Criteria and Internal Validation Parameters | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | $R^2$ | $R^2_{adj}$ | $Q^2_{LOO}$ | $RMSE_{Tr}$ | $S$ | $F$ | $CCC_{Tr}$ |
| | | | **Divison 1** | | | | | | |
| M1 | 6 | piPC08 CIC1 Chi_D/Dt L/Bw TDB09p Mor28s | 0.59 | 0.57 | 0.55 | 1.03 | 1.06 | 32.74 | 0.74 |
| **M2***[*] | **6** | **CIC2 Chi_D/Dt L/Bw TDB09p Mor28s piPC08** | **0.60** | **0.58** | **0.56** | **1.02** | **1.04** | **33.71** | **0.75** |
| | | | **Divison 2** | | | | | | |
| M3 | 6 | WiA_Dt P_VSA_MR_2 Mor31e B09[C-O] MLOGP2 HATS1m | 0.61 | 0.59 | 0.57 | 1.02 | 1.05 | 35.07 | 0.75 |
| | | | **Divison 3** | | | | | | |
| M4 | 3 | WiA_Dt Mor31e MLOGP2 | 0.56 | 0.55 | 0.53 | 1.07 | 1.09 | 45.04 | 0.72 |
| M5 | 6 | WiA_Dt P_VSA_MR_2 Mor31e R1m+ B09[C-O] MLOGP2 | 0.59 | 0.56 | 0.56 | 1.03 | 1.06 | 40.14 | 0.74 |

*Selected Model.

Threshold values: $R^2 > 0.6$, $Q^2_{LOO} > 0.5$

Table 4.4. The external parameters of the developed models for the mutagenic activity of textile dyes

| Model No | Number of Variables | Variables | Fitting Criteria and Internal Validation Parameters | | | | | | |
|----------|---------------------|-----------|-------------------|------|------|------|------|------|------|
| | | | $R^2_{Test}$ | $Q^2$-$_{F1}$ | $Q^2$-$_{F2}$ | $Q^2$-$_{F3}$ | $CCC_{Test}$ | $RMSE_{Test}$ | $MAE_{Test}$ |
| **Divison 1** | | | | | | | | | |
| M1 | 6 | piPC08 CIC1 Chi_D/Dt L/Bw TDB09p Mor28s | 0.74 | 0.75 | 0.74 | 0.83 | 0.86 | 0.65 | 0.55 |
| **M2**[*] | **6** | **CIC2 Chi_D/Dt L/Bw TDB09p Mor28s piPC08** | **0.76** | **0.75** | **0.74** | **0.84** | **0.85** | **0.65** | **0.53** |
| **Divison 2** | | | | | | | | | |
| M3 | 6 | WiA_Dt P_VSA_MR_2 Mor31e B09[C-O] MLOGP2 HATS1m | 0.77 | 0.76 | 0.76 | 0.86 | 0.87 | 0.60 | 0.45 |
| **Divison 3** | | | | | | | | | |
| M4 | 4 | WiA_Dt Mor31e MLOGP2 P_VSA_MR_2 | 0.74 | 0.73 | 0.73 | 0.85 | 0.85 | 0.64 | 0.52 |
| M5 | 5 | WiA_Dt Mor31e R1m+ B09[C-O] MLOGP2 | 0.62 | 0.75 | 0.74 | 0.85 | 0.86 | 0.62 | 0.45 |

*Selected Model.

Threshold values: $Q^2_{F1}$, $Q^2_{F2}$, $Q^2_{F3}$ > 0.70, $CCC_{Test}$ >

Table 4.5. The structural coverage of all QSAR models for external set chemicals and the results of Xternal Validation Plus tool (*MAE*-based criteria).

| Model Label | Number of compounds within the AD of Selected Model (out of 33) | Structural Coverage (%) | *MAE*-based Criteria |
|---|---|---|---|
| **Division 1** | | | |
| M1 | 22 | 66% | GOOD |
| **M2*** | **23** | **70%** | **GOOD** |
| **Division 2** | | | |
| M3 | 22 | 66% | GOOD |
| **Division 3** | | | |
| M4 | 19 | 58% | GOOD |
| M5 | 19 | 58% | GOOD |

**\*Selected Model**

Of the generated models, the 6-descriptor MLR model labelled as M2 was highlighted regarding its structural coverage (70%) for the external set chemicals with no mutagenicity data. Model M2 resulted from random division. The model includes 6 descriptors from DRAGON 7.0. software. The change in $R^2$ and $Q^2$ of the selected model was given in Figure 4.1, indicating that the addition of each descriptor leads to an increase in $R^2$ more than 0.02. In other words, there is no redundant descriptor in model M2.

Figure 4.1. The change in $R^2$ and $Q^2$ along with the increase in the number of variables.

Model M2 gives the following equation (Eq. 4.1) together with the descriptors involved and their regression coefficients and the 95% confidence intervals written in parenthesis.

**Log (mutagenicty)** $_{\text{(TA98+S9 rev/nmol)}}$ = -4.6629 (± 0.8191) + 0.9 (± 0.5918) **CIC2** + 1.0654 (± 0.2654)**Chi_D/Dt** – 0.0866 (±0.0558) **L/Bw** + 0.7161 (± 0.2631) **TDB09p** + 0.4636 (±0.2242) **Mor28s** – 0.1603 (±0.1405) **piPC08** (4.1)

$n_{\text{Tr}}$ = 142, $R^2$ = 0.60, $R^2_{\text{adj}}$=0.58, $RMSE_{\text{Tr}}$ = 1.02, $F$ = 33.71, $CCC_{\text{Tr}}$ = 0.75, $n_{\text{Test}}$ = 35, $R^2_{\text{Test}}$ = 0.74, $RMSE_{\text{Test}}$ = 0.65, $Q^2_{\text{F1}}$ = 0.75, $Q^2_{\text{F2}}$ = 0.74, $Q^2_{\text{F3}}$ = 0.84, $CCC_{\text{Test}}$ = 0.86

where the $n_{\text{Tr}}$ and $n_{\text{Test}}$ refer to the number of compounds in the training and test sets, respectively. The equations of all models except model M2 are given in Appendix C

Besides employed external validation, the prediction quality of the developed model was tested via *MAE*-based approach using Xternal Validation Plus tool developed by Roy et al. (2015). The output of this software for model M2 is given in Table 4.6, whereas the output of this software for the other models is given in Appendix D. Model M2 is regarded as **GOOD.** This criterion judges the performance of actual prediction quality of QSAR model based on the mean absolute error (*MAE*) showing that there was no systematic error in model M2.

Table 4.6. The results of Xternal Validation Plus tool for the selected model (M2).

| User Input File Info. | FileName | *MAE*-test M2.xlsx |
|---|---|---|
| Model biasness test | Systematic Error Result | Absent |
| | nPE / nNE | 0.7500 |
| | nNE / nPE | 1.3333 |
| | \|MPE / MNE\| | 0.9960 |
| | \|MNE / MPE\| | 1.0040 |
| | AAE - \|AE\| | 0.4559 |
| | $R^2$ (Residuals; serial correlation) | 0.0165 |
| | $R^2$ (Residuals and Yobs values) | 0.2080 |
| | $R^2_{Test}$(100% data) | 0.7431 |
| | $R_0^2{}_{Test}$(100% data) | 0.7385 |
| | $R_0'^2{}_{Test}$(100% data) | 0.7405 |
| Classical Metrics | $Q^2_{F1}$(100% data) | 0.7537 |
| (for 100% data) | $Q^2_{F2}$(100% data) | 0.7385 |
| | Scaled Avg.$R_m^2$(100% data) | 0.6455 |
| | Scaled Delta $R_m^2$(100% data) | 0.1930 |
| | $CCC$(100% data) | 0.8546 |
| | $R^2_{Test}$(95% data) | 0.8129 |
| | $R_0^2{}_{Test}$ (95% data) | 0.8128 |
| Classical Metric | $R_0'^2{}_{Test}$ (95% data) | 0.7426 |
| (after removing | $Q^2_{F1}$( 95% data) | 0.8198 |
| 5% data with | $Q^2_{F2}$ (95% data) | 0.8127 |
| high residuals) | Scaled Avg $R_m^2$(95% data) | 0.7389 |
| | Scaled Delta $R_m^2$(95% data) | 0.1317 |
| | $CCC$(95% data) | 0.8974 |
| | $RMSE_P$(100% data) | 0.6530 |
| Error-based metrics | SD(100% data) | 0.3826 |
| (for 100% data) | SE(100% data) | 0.0647 |
| | *MAE*(100% data) | 0.5330 |
| | $RMSE_P$(95% data) | 0.5542 |
| Error-based metric | SD(95% data) | 0.2940 |
| (after removing 5% data | SE(95% data) | 0.0512 |
| with high residuals) | *MAE*(95% data) | 0.4725 |
| | *MAE*+3*SD(95% data) | 1.3546 |

Table 4.6.  Continued.

| BASIC DATA STRUCTURE INFORMATION | | |
|---|---|---|
| | NCompTest | 35.0000 |
| Number of test set compounds, | Train range | 7.2900 |
| Range and Mean (train and test) | TrainY$_{Mean}$ | -0.2900 |
| | Test range | 5.4600 |
| | Test Y$_{Mean}$ | -0.6063 |
| Distribution of observed | %Y(+/-0.5)TestMean | 31.4286 |
| response values of Test set | %Y(+/-1.0)TestMean | 62.8571 |
| around Test mean(in %) | %Y(+/-1.5)TestMean | 80.0000 |
| | %Y(+/-2.0)TestMean | 85.7143 |
| Distribution of observed | %Y(+/-0.5)TrainMean | 37.1429 |
| response values of Test set | %Y(+/-1.0)TrainMean | 54.2857 |
| around Train mean (in %) | %Y(+/-1.5)TrainMean | 77.1429 |
| | %Y(+/-2.0)TrainMean | 85.7143 |
| | %NComp>(0.1*TR) | 25.7143 |
| Distribution of prediction | %NComp>(0.15*TR) | 8.5714 |
| errors (in %) | %NComp>(0.2*TR) | 2.8571 |
| | %NComp>(0.25*TR) | 0.0000 |
| | (0.1*Training Set Range) | 0.7290 |
| Threshold values utilized | (0.15*Training Set Range) | 1.0935 |
| to judge the model predictions | (0.2*Training Set Range) | 1.4580 |
| | (0.25*Training Set Range) | 1.8225 |
| RESULT (*MAE*-based criteria applied on 95% data) | Prediction Quality | **GOOD** |

From a general notation, an error of 10% of the training set range should be acceptable while an error value more than 20% of the training set should be a very high error. Thus, the criteria for good predictions should be the following:

*MAE* ≤ 0.1 x training set range and *MAE*+ 3σ ≤ 0.2 x training set range.
**0.4708 ≤ 0.1 x 7.29  and  1.3283 ≤ 0.2 x 7.29**

where, the σ value refers to the standard deviation of the absolute error values for the test set data. Considering a normal distribution pattern, mean ± 3σ covers 99.7% of the data points.

Predicted vs. observed mutagenicity values were plotted in Figure 4.2, together with y=x line. The training and test set compounds distributed homogeneously around the line (Figure 4.2.)

Figure 4.2. The plot of calculated/predicted vs. observed values of mutagenicity for the training/test set compounds by model equation (Eq.4.1), with yellow labeled training set compounds and blue labeled test set compounds.

The proposed model (M2) includes 6 descriptors from DRAGON 7.0 (Kode, 2017). The description of these 6 parameters is listed in Table 4.7. The description of descriptors appearing in the other four models are given in Table C1.

Regardless of their sign, the importance of 6 descriptors could be written and explained as below based on the magnitude of standardized coefficients:

Chi_D/Dt > CIC2 > TDB09p > Mor28s > piPC08 > L/Bw

Table 4.7. Descriptors appeared in model M2 and their descriptions.

| Abbreviation of Descriptor | Description | Block |
|---|---|---|
| CIC2 | Complementary Information Content Index (neighborhood symmetry of 2-order) | Information Indices |
| Chi_D/Dt | Randic-like index from distance/detour matrix | 2D matrix-based descriptors |
| L/Bw | The length-to-breadth ratio by WHIM | Geometrical descriptors |
| TDB09p | 3D Topological distance-based descriptors - lag 9 weighted by polarizability | 3D autocorrelations |
| Mor28s | signal 28 / weighted by I-state | 3D-MoRSE descriptor* |
| piPC08 | molecular multiple path count of order 8 | Walk and path counts |

*MoRSE: Molecule Representation of Structure based on Electron diffraction.

The most important descriptor affecting the mutagenicity of chemicals is, Chi_D/Dt, Randic-like index from distance/detour matrix. Topochemical indices usually are two dimensional (2D) indices and encode information pertaining to both molecular topology and chemical nature of atoms and bonds in a molecule (Basak and Gute, 1997). Distance-based topological indices employ the distance matrix or the detour matrix to characterise molecular graphs. The distance matrix is based on topological distance, which is the number of edges in the shortest path between vertices $v_i$ and $v_j$, whereas the detour matrix is based on the number of edges in the longest path between vertices $v_i$ and $v_j$ in a molecular graph G (Todeschini and Consonni, 2000). All these information reveal that as the distance between the atoms in 2D graph of a chemical increseas, mutagenicity increases.

The second important descriptor is CIC2 from information content block. The CIC2 index provides information about the abundance of rings in a molecule (Todeschini and Consonni, 2000). IC2 is the second-grade classification sphere near the atomic radius (Song et al., 2015). Both CIC2 and Chi_D/Dt values are positively correlated to the mutagenicity of chemicals in the dataset. Similar to our finding, Song et al. (2015) indicated that CIC2 index is positively correlated to the toxicity ($LC_{50}$) of nitrobenzene to *Paramecium caudatum*.

TDB09p and Mor28s are 3D descriptors and they are also positively correlated to mutagenicity. TDB09p which is the 3D topological distance-based descriptors -lag 9 weighted by polarizability has positive coefficient in the model. It is likely that chemicals with more negatively charged/ionizable group, lipophilic group and carbon-sulfur atom and with large 3D topological distance tended to have

high mutagenic activity. TDB09 descriptor appeared in a QSAR model developed for predicting the binding affinity of endocrine disrupting chemicals to eight fish estrogen receptor (He et al., 2018). On the contrary, it has a negative coefficient in their QSAR model indicating that chemicals with more negatively/charged/ionizable group, lipophilic group and carbon-sulfur atom and with large 3D topological distance tended not to bind with estrogen receptor. This difference can be attributed to the different targets of chemicals in the two dataset.

The other 3D descriptor, Mor28s, is related to molecule representation of structures based on electron diffraction (MoRSE) weighted by electric state/I-state. It is based on the three-dimensional structure of molecules by a certain number of values and contains some molecular codes obtained by mathematical transformations used in electron diffraction (Schuur et al., 1996). In a modelling study of pharmaceuticals to fish, Tugcu et al. (2012) have used 3D-MoRSE descriptors and reported that 3D-MoRSE descriptors are related to fish toxicity because of its correlation with hydrophobicity. In addition, Caballero and Fernandez (2011) have also used a 3D-MoRSE descriptor for antifungal activity modeling.

The two descriptors L/Bw from WHIM descriptor group and piPC08 from walk and path counts group have inverse relationships with the mutagenicity. The negative coefficients of geometric descriptor (L/Bw) and molecular multiple path count of order 8 (piPC08) signify that as the value of these descriptors increases, mutagenicity decreases. The descriptor, piPC08, has a negative coefficient in a QSAR model equation in the study of Kusic et al. (2008) in which they modeled the rate constants for radical degradation of aromatic pollutants in water matrix. Taken together, these data suggest that descriptors relevant to the 2D and 3D topology, abundance of rings in a molecule, 3D geometry, polarizability which is reflected by the more negatively charged/ionizable group, lipophilic group and carbon-sulfur atom and with large 3D topological distance, and multiple path count of order 8 affected the mutagenic activity of the studied textile dyes.

## 4.2. Applicability Domain of the Selected Model

The ADs of linear models were defined by the boundaries of the descriptor and the response range (Table 4.8). Williams plot (Figure 4.3) shows the hat values of chemicals and their prediction accuracy. Errors are represented by standardized residuals. The vertical reference line refers to the critical hat value ($h$*=0. 0.148) and the horizontal reference lines are ±3σ, the cut-off values for the response outliers. While model M2 has no response outlier, it has only one structural outlier. Hat values of all chemicals were lower than the critical hat value ($h$* = 0.148), except 2-bromo-7-

aminofluorene Although this compound (2-bromo-7-aminofluorene) was out of the descriptor range, their predicted mutagenic activity value is good enough. Therefore we included it in our dataset, to increase the AD of the proposed model.



Figure 4.3.  Williams plot for Model M2 with training set in yellow and test set in blue. ID number, 70 refers to 2-bromo-7-aminofluorene, respectively.

Table 4.8.  Ranges of descriptors used in the model.

| Variable | Minimum value | Maximum value |
| --- | --- | --- |
| Chi_D/Dt | 0.083 | 0.277 |
| CIC2 | 0.450 | 2.236 |
| TDB09p | 0.000 | 4.612 |
| Mor28s | -1.096 | 3.792 |
| piPC08 | 0.000 | 8.230 |
| L/Bw | 1.210 | 23.36 |

Status of chemicals in the dataset, experimental and predicted mutagenicity values from model M2, descriptor and  hat values were given in Table 4.9

Table 4.9. Chemicals that are used to model mutagenicity, their experimental and predicted values, descriptor and hat values.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 4'-NEt2-3-OMe-AAB | Training | -1.77 | -2.16 | 0.05 | 1.53 | 0.15 | 10.78 | 4.18 | 0.48 | 6.28 |
| 2 | 2-OMe-AAB | Training | 3.31 | 3.00 | 0.07 | 1.41 | 0.17 | 6.15 | 2.92 | 1.04 | 6.00 |
| 3 | 4'-OH-AAB | Prediction | 0.41 | 0.39 | 0.06 | 1.56 | 0.18 | 15.15 | 3.37 | 2.00 | 5.96 |
| 4 | 3'-Me-4'-OH-AAB | Prediction | 0.60 | 0.99 | 0.04 | 1.24 | 0.17 | 11.04 | 3.29 | 1.47 | 6.03 |
| 5 | 4'-OH-2',3-diMe-AAB(4'-OH-OAT) | Training | -2.52 | -1.86 | 0.05 | 1.18 | 0.17 | 8.23 | 3.24 | 1.71 | 6.13 |
| 6 | AAB | Training | 0.87 | 1.24 | 0.03 | 1.90 | 0.18 | 22.72 | 3.27 | 0.89 | 5.84 |
| 7 | 3'-Me-AAB | Prediction | 3.00 | 1.02 | 0.06 | 1.43 | 0.17 | 15.64 | 3.26 | 0.97 | 5.93 |
| 8 | 3-OMe-4'-N(CH2CH2OH)2-AAB | Prediction | -1.52 | -1.16 | 0.04 | 1.39 | 0.15 | 11.80 | 4.11 | 1.84 | 6.34 |
| 9 | 3'CH2OH-AAB | Training | 3.23 | 2.90 | 0.06 | 1.29 | 0.17 | 8.35 | 3.31 | 1.78 | 5.99 |
| 10 | 3-OH-AAB | Training | 0.48 | 1.39 | 0.03 | 1.46 | 0.18 | 10.69 | 3.17 | 0.53 | 5.93 |
| 11 | 3-OCH2CH2OH-AAB | Training | -1.51 | -0.91 | 0.06 | 1.40 | 0.16 | 8.35 | 3.89 | 1.40 | 6.10 |
| 12 | 2'-CH2OH-3-Me-AAB | Prediction | -0.67 | -0.59 | 0.05 | 1.04 | 0.16 | 5.26 | 3.12 | 1.57 | 6.08 |
| 13 | 4-OMe-AAB | Training | 3.16 | 3.08 | 0.07 | 1.56 | 0.17 | 16.11 | 3.83 | 1.73 | 5.99 |
| 14 | 2',3-diMe-AAB | Training | -0.48 | 1.02 | 0.04 | 1.41 | 0.17 | 7.42 | 3.19 | 1.41 | 6.04 |
| 15 | 3-OBu-AAB | Prediction | -3.00 | -1.48 | 0.05 | 1.39 | 0.15 | 6.96 | 3.16 | 1.04 | 6.12 |
| 16 | 3-OEt-AAB | Training | 1.12 | 1.38 | 0.02 | 1.40 | 0.16 | 5.33 | 3.70 | 0.92 | 6.07 |
| 17 | 3-OPr-AAB | Training | -1.17 | -0.97 | 0.06 | 1.34 | 0.16 | 3.56 | 3.46 | 0.88 | 6.10 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|-----------------------------------|---------------------------------|------------------------|------|----------|------|--------|--------|--------|
| 18 | 3-OMe-AAB | Training | 1.83 | 3.97 | 0.12 | 1.48 | 0.17 | 6.81 | 3.52 | 1.03 | 6.01 |
| 19 | 3'-Me-4'-OH-MAB | Training | 3.97 | 1.23 | 0.02 | 1.21 | 0.17 | 12.64 | 3.62 | 1.40 | 6.06 |
| 20 | 3'-COOH-MAB | Training | -3.32 | -1.74 | 0.04 | 1.24 | 0.16 | 12.39 | 3.91 | 1.63 | 6.13 |
| 21 | 4'-OH-MAB | Training | 3.35 | 1.94 | 0.07 | 1.49 | 0.17 | 17.65 | 3.85 | 1.76 | 5.99 |
| 22 | MAB | Training | 2.46 | 2.00 | 0.04 | 1.79 | 0.18 | 14.02 | 3.78 | 1.30 | 5.88 |
| 23 | 4'-Me-MAB | Training | 3.50 | 3.38 | 0.08 | 1.55 | 0.17 | 16.07 | 4.28 | 1.04 | 5.99 |
| 24 | 3'-Me-MAB | Training | 1.18 | 0.90 | 0.07 | 1.38 | 0.17 | 12.15 | 3.58 | 0.87 | 5.96 |
| 25 | 3'-CH2OH-MAB | Training | -1.04 | 0.27 | 0.06 | 1.26 | 0.16 | 10.08 | 3.64 | 1.57 | 6.02 |
| 26 | 3'-Me-4'-OH-DAB | Training | -1.77 | -0.21 | 0.07 | 1.46 | 0.16 | 13.29 | 3.89 | 1.19 | 6.09 |
| 27 | *p*(dimethylamino)azobenzene | Training | 0.89 | 1.19 | 0.04 | 2.03 | 0.17 | 14.82 | 4.18 | 1.17 | 5.92 |
| 28 | 3'-COOH-DAB | Training | 0.38 | -0.78 | 0.04 | 1.49 | 0.16 | 12.25 | 4.18 | 1.40 | 6.16 |
| 29 | 2-Me-DAB | Prediction | -3.14 | -1.45 | 0.09 | 1.65 | 0.17 | 10.20 | 3.59 | 0.84 | 6.01 |
| 30 | 3'-Me-DAB | Training | 2.88 | 2.13 | 0.07 | 1.63 | 0.17 | 13.00 | 3.86 | 0.65 | 6.00 |
| 31 | 3'-CHO-DAB | Training | -2.70 | -2.51 | 0.12 | 1.54 | 0.16 | 10.71 | 4.14 | 1.54 | 6.11 |
| 32 | 3'-CH2-OAC-DAB | Training | -3.00 | -2.24 | 0.05 | 1.44 | 0.14 | 6.81 | 3.90 | 1.28 | 6.20 |
| 33 | 3'-CH2OH-DAB | Training | -0.01 | 0.80 | 0.04 | 1.51 | 0.16 | 9.83 | 3.91 | 1.37 | 6.06 |
| 34 | 3'-Me-AAB-N-Ac | Training | -1.60 | -1.79 | 0.05 | 1.48 | 0.16 | 15.58 | 3.65 | 0.84 | 6.16 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 35 | 3'-Me-4'-OH-AAB-N-Ac | Training | -1.60 | -0.53 | 0.04 | 1.31 | 0.16 | 14.56 | 3.69 | 1.17 | 6.25 |
| 36 | N-OH-2-OMe-AAB | Training | -1.42 | 0.47 | 0.05 | 1.30 | 0.17 | 7.30 | 2.97 | 0.85 | 6.04 |
| 37 | 3'-Me-MAB-N-Ac | Training | -2.05 | -2.42 | 0.04 | 1.49 | 0.16 | 16.06 | 3.92 | 0.67 | 6.19 |
| 38 | N-OH-MAB | Prediction | 1.13 | 0.51 | 0.04 | 1.74 | 0.17 | 14.67 | 3.85 | 1.37 | 5.92 |
| 39 | N-OH-3'-Me-MAB | Prediction | 3.80 | 2.74 | 0.06 | 1.34 | 0.17 | 13.39 | 3.60 | 0.95 | 6.00 |
| 40 | N-OH-AAB | Training | 2.38 | 1.75 | 0.03 | 1.75 | 0.18 | 13.84 | 3.34 | 1.25 | 5.88 |
| 41 | N-OH-4'-Me-MAB | Prediction | 1.43 | 2.99 | 0.07 | 1.49 | 0.17 | 16.36 | 3.97 | 0.81 | 5.99 |
| 42 | N-OH-3-OMe-AAB | Training | 0.04 | 0.92 | 0.03 | 1.36 | 0.17 | 11.37 | 3.81 | 1.14 | 6.06 |
| 43 | 2-aminobiphenyl | Training | 2.62 | 1.31 | 0.03 | 1.95 | 0.19 | 4.50 | 1.34 | 0.44 | 5.77 |
| 44 | 4-aminobiphenyl | Training | 1.18 | 1.24 | 0.03 | 1.97 | 0.20 | 8.78 | 2.06 | 0.16 | 5.77 |
| 45 | 2-2'-diaminobiphenyl | Training | 2.98 | 1.95 | 0.04 | 1.92 | 0.18 | 3.64 | 1.35 | 0.91 | 5.90 |
| 46 | 2,4'-diaminobiphenyl | Training | 0.43 | 0.43 | 0.05 | 1.75 | 0.19 | 6.32 | 1.67 | 0.98 | 5.90 |
| 47 | 3,3'-diamnobiphenyl | Training | 1.93 | 1.12 | 0.04 | 1.77 | 0.19 | 5.02 | 1.26 | 1.10 | 5.90 |
| 48 | 3,4'-diaminobiphenyl | Training | 2.62 | 1.68 | 0.16 | 1.83 | 0.19 | 7.21 | 2.29 | 0.66 | 5.94 |
| 49 | 2-amino-3'-nitrobiphenyl | Training | -0.60 | -0.14 | 0.04 | 1.38 | 0.17 | 5.21 | 1.53 | 1.13 | 6.07 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|----------------------------------|--------------------------------|------------------------|------|----------|------|--------|--------|--------|
| 50 | 2-amino-4'nitrobiphenyl | Prediction | -2.40 | -1.40 | 0.05 | 1.44 | 0.18 | 8.48 | 2.44 | 1.98 | 6.19 |
| 51 | 3-amino-2'nitrobiphenyl | Training | -2.30 | -2.03 | 0.03 | 1.41 | 0.17 | 2.61 | 1.31 | 1.93 | 6.14 |
| 52 | 3-amino-3'nitrobiphenyl | Training | -2.15 | -0.24 | 0.06 | 1.46 | 0.18 | 4.78 | 1.42 | 1.89 | 6.23 |
| 53 | 3-amino-4'nitrobiphenyl | Training | -0.54 | -1.30 | 0.13 | 1.52 | 0.18 | 10.38 | 2.63 | 2.23 | 6.22 |
| 54 | 4-amino-2'nitrobiphenyl | Training | -2.22 | -1.61 | 0.09 | 1.52 | 0.17 | 3.21 | 1.74 | 1.44 | 6.21 |
| 55 | 4-amino-3'nitrobiphenyl | Training | -1.32 | -1.11 | 0.08 | 1.52 | 0.18 | 6.48 | 2.40 | 1.64 | 6.22 |
| 56 | 4-amino-4'nitrobiphenyl | Training | -2.22 | -1.55 | 0.05 | 1.77 | 0.18 | 14.39 | 2.37 | 1.58 | 6.19 |
| 57 | 1-aminoanthracene | Prediction | -2.10 | -2.15 | 0.06 | 1.80 | 0.11 | 3.59 | 1.31 | 0.48 | 7.42 |
| 58 | 2-aminoanthracene | Training | -2.52 | -1.97 | 0.05 | 1.82 | 0.11 | 5.71 | 2.04 | 0.10 | 7.45 |
| 59 | 9-aminoanthracene | Prediction | -3.00 | -2.34 | 0.08 | 1.80 | 0.11 | 3.59 | 1.31 | 0.48 | 7.42 |
| 60 | 1-aminophenanthrene | Training | -1.42 | -1.40 | 0.04 | 1.88 | 0.11 | 3.04 | 1.31 | 0.96 | 7.40 |
| 61 | 2-aminophenanthrene | Prediction | -3.00 | -2.12 | 0.04 | 1.85 | 0.11 | 3.97 | 2.07 | 0.79 | 7.41 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|----------------------------------|-------------------------------|------------------------|------|----------|------|--------|--------|--------|
| 62 | 3-aminophenanthrene | Training | -0.80 | -1.82 | 0.04 | 1.85 | 0.11 | 2.83 | 1.32 | 0.71 | 7.41 |
| 63 | 9-aminophenanthrene | Training | -0.87 | -2.03 | 0.03 | 1.87 | 0.11 | 2.01 | 1.35 | 1.17 | 7.38 |
| 64 | 1-aminofluorene | Training | -3.00 | -1.45 | 0.02 | 1.30 | 0.12 | 3.25 | 0.00 | 0.97 | 6.64 |
| 65 | 2-aminofluorene | Prediction | -0.49 | -1.35 | 0.04 | 1.27 | 0.12 | 5.06 | 1.40 | 0.95 | 6.65 |
| 66 | 3-aminofluorene | Prediction | -0.77 | -1.47 | 0.04 | 1.27 | 0.12 | 3.85 | 1.35 | 0.91 | 6.66 |
| 67 | 4-aminofluorene | Training | -0.21 | -1.67 | 0.03 | 1.30 | 0.12 | 2.58 | 0.00 | 0.97 | 6.64 |
| 68 | 2,7-diaminofluorene | Training | -1.80 | -1.76 | 0.02 | 1.46 | 0.12 | 6.58 | 1.89 | 1.39 | 6.76 |
| 69 | 2-amino-7-nitrofluorene | Training | -0.69 | -0.98 | 0.05 | 1.17 | 0.11 | 8.85 | 2.10 | 2.62 | 7.00 |
| 70 | 2-bromo-7-aminofluorene | Training | -2.40 | -1.11 | 0.08 | 1.02 | 0.12 | 15.46 | 4.20 | 1.01 | 6.76 |
| 71 | 2-hydroxy-7-aminofluorene | Training | 0.35 | -0.87 | 0.07 | 0.98 | 0.12 | 6.65 | 1.66 | 0.59 | 6.76 |
| 72 | 2-amino-7-acetamidofluorene | Training | -1.43 | -1.76 | 0.02 | 1.07 | 0.11 | 10.52 | 3.82 | 0.79 | 6.97 |
| 73 | 1-aminopyrene | Training | -0.40 | -0.17 | 0.05 | 2.03 | 0.10 | 2.02 | 1.31 | 1.25 | 8.21 |
| 74 | 2-aminopyrene | Training | -0.47 | -0.59 | 0.03 | 2.07 | 0.10 | 2.51 | 2.06 | 1.13 | 8.21 |
| 75 | 4-aminopyrene | Training | -0.59 | -0.12 | 0.04 | 1.96 | 0.10 | 1.41 | 1.34 | 1.34 | 8.20 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|------------------------------|-------------------------------|------------------------|------|----------|------|--------|--------|--------|
| 76 | 1-aminofluoranthene | Training | -0.24 | -0.27 | 0.03 | 1.85 | 0.10 | 1.50 | 0.00 | 1.09 | 8.22 |
| 77 | 2-aminofluoranthene | Training | 0.47 | 0.07 | 0.03 | 1.90 | 0.10 | 1.82 | 1.36 | 1.17 | 8.23 |
| 78 | 3-aminofluoranthene | Training | 0.64 | -0.66 | 0.08 | 1.96 | 0.10 | 2.53 | 1.40 | 1.30 | 8.23 |
| 79 | 7-aminofluoranthene | Training | -0.18 | -0.32 | 0.06 | 1.85 | 0.10 | 1.93 | 0.00 | 1.62 | 8.22 |
| 80 | 8-aminofluoranthene | Training | 0.13 | -0.57 | 0.05 | 1.86 | 0.10 | 2.71 | 1.39 | 1.16 | 8.23 |
| 81 | 6-aminochrysene | Training | -0.64 | -0.73 | 0.06 | 2.24 | 0.08 | 3.47 | 2.89 | 1.38 | 8.18 |
| 82 | 3-aminoquinoline | Training | 0.55 | 1.12 | 0.03 | 0.82 | 0.16 | 3.35 | 0.00 | -0.06 | 6.06 |
| 83 | 5-aminoquinoline | Training | -2.67 | -1.39 | 0.09 | 0.72 | 0.16 | 1.71 | 0.00 | 0.51 | 6.02 |
| 84 | 6-aminoquinoline | Training | 0.20 | 0.81 | 0.03 | 0.72 | 0.16 | 3.37 | 0.00 | 0.31 | 6.06 |
| 85 | 8-aminoquinoline | Prediction | -2.00 | -1.05 | 0.08 | 0.82 | 0.16 | 1.68 | 0.00 | 0.14 | 6.02 |
| 86 | 1-aminocarbazole | Prediction | -2.15 | -1.58 | 0.05 | 1.27 | 0.12 | 3.22 | 0.00 | 0.95 | 7.37 |
| 87 | 2-aminocarbazole | Training | -0.38 | -1.60 | 0.04 | 1.24 | 0.12 | 5.11 | 1.40 | 1.00 | 7.38 |
| 88 | 3-aminocarbazole | Training | -1.29 | -1.41 | 0.04 | 1.24 | 0.12 | 3.96 | 1.36 | 0.88 | 7.39 |
| 89 | 4-aminocarbazole | Training | -0.87 | -1.31 | 0.06 | 1.27 | 0.12 | 2.58 | 0.00 | 1.20 | 7.37 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 90 | 1-aminophenazine | Training | -2.70 | -1.60 | 0.02 | 1.58 | 0.11 | 3.43 | 1.29 | -0.01 | 7.42 |
| 91 | 2-aminophenazine | Training | -1.40 | -2.28 | 0.07 | 1.60 | 0.11 | 5.50 | 2.01 | 0.14 | 7.45 |
| 92 | 1,6-diaminophenazine | Training | -0.39 | -0.73 | 0.07 | 1.52 | 0.11 | 3.02 | 1.24 | 0.56 | 7.50 |
| 93 | 1,7-diaminophenazine | Prediction | 0.81 | 0.33 | 0.03 | 1.40 | 0.11 | 4.42 | 2.20 | 0.66 | 7.55 |
| 94 | 1,9-diaminophenazine | Training | -0.04 | 0.41 | 0.08 | 1.52 | 0.11 | 2.77 | 1.27 | 0.70 | 7.51 |
| 95 | 2,7-diaminophenazine | Training | 0.15 | -0.51 | 0.02 | 1.67 | 0.11 | 6.74 | 2.15 | 0.88 | 7.56 |
| 96 | 2,8-diaminophenazine | Training | 0.01 | -0.53 | 0.02 | 1.67 | 0.11 | 6.19 | 2.38 | 0.71 | 7.56 |
| 97 | 1-naphthylamine | Training | -1.14 | -1.12 | 0.08 | 1.56 | 0.16 | 1.73 | 0.00 | 0.76 | 6.02 |
| 98 | 2-naphthylamine | Prediction | 0.75 | 1.19 | 0.03 | 1.52 | 0.16 | 3.41 | 0.00 | 0.39 | 6.06 |
| 99 | 1-amino-4-nitronaphthalene | Training | 0.32 | -1.85 | 0.03 | 1.25 | 0.15 | 1.70 | 0.00 | 2.83 | 6.47 |
| 100 | 2-amino-1-nitronaphthalene | Training | -0.02 | -2.26 | 0.04 | 1.19 | 0.15 | 1.84 | 0.00 | 1.37 | 6.54 |
| 101 | aniline | Prediction | -0.16 | -0.67 | 0.01 | 1.46 | 0.28 | 2.14 | 0.00 | 0.34 | 0.00 |
| 102 | *p*-bromoaniline | Training | 0.13 | 0.22 | 0.03 | 1.00 | 0.27 | 14.02 | 0.00 | 0.45 | 0.00 |
| 103 | *o*-chloroaniline | Training | 0.43 | -0.03 | 0.02 | 0.86 | 0.27 | 2.20 | 0.00 | 0.64 | 0.00 |
| 104 | *p*-chloroaniline | Training | 0.30 | -0.20 | 0.04 | 1.00 | 0.27 | 6.85 | 0.00 | 0.52 | 0.00 |
| 105 | *p*-fluoroaniline | Training | -0.22 | -0.02 | 0.03 | 1.00 | 0.27 | 4.07 | 0.00 | 0.26 | 0.00 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 106 | *p*-ethoxyaniline | Training | -0.96 | -0.62 | 0.02 | 0.99 | 0.25 | 5.53 | 1.07 | -0.07 | 2.78 |
| 107 | 4-phenoxyaniline | Training | 0.69 | -0.18 | 0.03 | 1.89 | 0.19 | 9.10 | 2.33 | -0.07 | 5.56 |
| 108 | 4-hydroxyaniline | Prediction | -0.62 | -0.87 | 0.04 | 0.93 | 0.27 | 3.70 | 0.00 | 0.22 | 0.00 |
| 109 | 4-cyclohexylaniline | Training | 0.20 | -0.32 | 0.03 | 1.80 | 0.20 | 12.04 | 1.79 | -0.21 | 4.47 |
| 110 | 2,4-dinitroaniline | Training | -0.62 | -0.30 | 0.02 | 1.04 | 0.23 | 2.46 | 0.00 | 3.79 | 4.72 |
| 111 | 2,4-dimethylaniline | Training | -2.23 | -0.48 | 0.02 | 1.21 | 0.26 | 2.23 | 0.00 | 0.06 | 0.00 |
| 112 | 2,5-dimethylaniline | Training | -0.69 | -1.03 | 0.14 | 1.31 | 0.26 | 2.29 | 0.00 | 0.20 | 0.00 |
| 113 | 2,4-difluoroaniline | Training | -1.15 | -0.51 | 0.03 | 0.63 | 0.26 | 2.14 | 0.00 | 0.82 | 0.00 |
| 114 | 2-methoxy-5-methylaniline | Prediction | 2.71 | 0.18 | 0.04 | 0.78 | 0.25 | 2.96 | 0.00 | 0.23 | 2.15 |
| 115 | 3-methoxy-4-methylaniline | Training | -0.45 | -0.32 | 0.02 | 0.78 | 0.25 | 1.62 | 0.00 | -0.24 | 2.15 |
| 116 | 4-methoxy-2-methylaniline | Prediction | -1.28 | -0.19 | 0.04 | 0.78 | 0.25 | 3.38 | 0.00 | -0.01 | 0.00 |
| 117 | 4-chloro-2-nitroaniline | Training | 0.01 | -0.27 | 0.03 | 0.55 | 0.24 | 2.28 | 0.00 | 2.94 | 3.45 |
| 118 | 2,4,5-trimethylaniline | Training | -0.95 | -0.22 | 0.03 | 1.71 | 0.25 | 2.09 | 0.00 | 0.16 | 0.00 |
| 119 | 2-bromo-4,6-dinitroaniline | Training | -0.89 | -0.99 | 0.03 | 0.78 | 0.22 | 1.46 | 0.00 | 3.75 | 4.72 |
| 120 | 4,4'-ethylenebis(aniline) | Training | -1.30 | -0.48 | 0.04 | 2.00 | 0.18 | 15.93 | 3.61 | 0.53 | 5.27 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 121 | 4,4'-methylenedianiline | Training | 1.04 | -0.81 | 0.06 | 1.93 | 0.18 | 9.21 | 2.39 | 0.57 | 5.67 |
| 122 | 4,4'-methylenebis (O-Ethylaniline) | Training | -0.42 | 0.35 | 0.03 | 0.94 | 0.25 | 1.96 | 0.00 | 0.23 | 2.15 |
| 123 | 4,4'-methylenebis(o-fluoroaniline) | Training | 0.70 | -0.41 | 0.02 | 0.86 | 0.27 | 1.42 | 0.00 | 0.65 | 0.00 |
| 124 | 4,4'-methylenebis(o-isopropylaniline) | Training | -1.05 | -0.78 | 0.03 | 1.28 | 0.25 | 2.15 | 0.00 | -0.15 | 2.78 |
| 125 | 1,2-phenylenediamine | Training | -0.74 | 0.10 | 0.03 | 1.50 | 0.27 | 1.44 | 0.00 | 1.03 | 0.00 |
| 126 | 1,3-benzenediamine | Prediction | -0.91 | -0.19 | 0.03 | 1.38 | 0.27 | 2.03 | 0.00 | 0.88 | 0.00 |
| 127 | 1,4-benzenediamine | Training | -0.85 | -0.32 | 0.06 | 1.75 | 0.27 | 3.69 | 0.00 | 0.98 | 0.00 |
| 128 | 2-methyl-1,4-phenylenediamine | Training | -0.35 | -0.51 | 0.02 | 1.13 | 0.26 | 2.25 | 0.00 | 1.07 | 0.00 |
| 129 | 2-nitro-1,4-phenylenediamine | Training | 2.28 | -0.21 | 0.02 | 1.04 | 0.24 | 2.05 | 0.00 | 3.25 | 3.45 |
| 130 | 2-methoxy-1,4-phenylenediamine | Training | -0.22 | 0.17 | 0.03 | 1.08 | 0.25 | 1.59 | 0.00 | 0.58 | 2.15 |
| 131 | 2-ethoxy-1,4-phenylenediamine | Training | -0.19 | -0.01 | 0.03 | 1.02 | 0.24 | 1.57 | 0.00 | 0.45 | 3.36 |
| 132 | 2-propoxy-1,4-phenylenediamine | Training | -1.49 | -0.53 | 0.05 | 0.98 | 0.22 | 2.11 | 1.05 | 0.59 | 3.66 |

Table 4.9. Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|----------|----------|----------|------|----------|------|--------|--------|--------|
| 133 | 2-butyloxy-1,4-phenylenediamine | Training | -0.92 | -0.55 | 0.03 | 1.09 | 0.21 | 4.47 | 1.51 | 0.85 | 3.87 |
| 134 | 4-chloro-1,2-phenylenediamine | Prediction | -0.85 | 0.40 | 0.05 | 1.05 | 0.26 | 4.11 | 0.00 | 1.18 | 0.00 |
| 135 | 4-chloro-1,3-phenylenediamine | Training | -0.55 | -0.16 | 0.04 | 0.92 | 0.26 | 3.43 | 0.00 | 1.05 | 0.00 |
| 136 | 4-nitro-1,2-phenylenediamine | Training | -0.66 | -0.13 | 0.02 | 1.15 | 0.25 | 3.73 | 0.00 | 2.46 | 0.00 |
| 137 | 4-nitro-1,3-phenylenediamine | Training | 0.05 | -0.56 | 0.04 | 1.04 | 0.25 | 3.10 | 0.00 | 3.20 | 3.45 |
| 138 | *N*-acetyl-1,4-phenylenediamine | Prediction | -0.14 | -0.62 | 0.05 | 1.08 | 0.24 | 6.28 | 1.30 | 0.70 | 3.84 |
| 139 | 1,4-*N*Ac-phenylenediamine | Prediction | -1.52 | -0.31 | 0.05 | 1.08 | 0.24 | 6.28 | 1.30 | 0.70 | 3.84 |
| 140 | 2,6-dichloro-1,4-phenylenediamine | Training | -0.70 | 0.11 | 0.03 | 1.13 | 0.25 | 1.21 | 0.00 | 1.41 | 0.00 |
| 141 | *N,N*-diethyl-1,4-phenylenediamine | Training | 0.00 | -0.67 | 0.02 | 1.63 | 0.25 | 3.75 | 1.11 | -0.17 | 3.45 |
| 142 | *N,N*-Dimethyl-1,4-phenylenediamine | Training | 1.14 | 0.18 | 0.04 | 1.61 | 0.26 | 4.50 | 0.00 | 0.67 | 0.00 |
| 143 | *N*-methyl-1,4-phenylenediamine | Training | 1.28 | 0.03 | 0.05 | 1.20 | 0.26 | 4.95 | 0.00 | 0.70 | 0.00 |
| 144 | 2,4-diaminoethylbenzene | Training | -1.06 | -0.84 | 0.03 | 0.98 | 0.25 | 3.06 | 0.00 | 0.68 | 2.15 |

Table 4.9.  Continued.

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 145 | 2,4-diaminotoluene | Training | 1.89 | 0.10 | 0.03 | 1.03 | 0.26 | 2.31 | 0.00 | 0.76 | 0.00 |
| 146 | 3,4-diaminotoluene | Training | 0.36 | -0.15 | 0.04 | 1.13 | 0.26 | 2.32 | 0.00 | 0.58 | 0.00 |
| 147 | 3-amino-α,α,α-trifluorotoluene | Training | -0.30 | -0.13 | 0.03 | 0.87 | 0.26 | 3.95 | 0.00 | 1.11 | 0.00 |
| 148 | 2,4-diamino-isopropylbenzene | Prediction | -0.28 | -0.81 | 0.04 | 1.29 | 0.25 | 3.16 | 0.00 | 0.31 | 2.78 |
| 149 | 2,4-diamino-n-butylbenzene | Training | -1.30 | -0.66 | 0.04 | 1.16 | 0.23 | 7.17 | 1.85 | 0.49 | 3.88 |
| 150 | 2-amino-4-chlorophenol | Prediction | -0.41 | 0.00 | 0.04 | 0.45 | 0.26 | 3.91 | 0.00 | 0.23 | 0.00 |
| 151 | 2-amino-4-methylphenol | Training | -2.00 | -0.28 | 0.02 | 0.64 | 0.26 | 2.41 | 0.00 | -0.04 | 0.00 |
| 152 | 2-amino-5-nitrophenol | Training | -0.96 | -0.29 | 0.02 | 0.63 | 0.25 | 3.66 | 0.00 | 1.14 | 0.00 |
| 153 | 3-amino-6-methyl phenol | Training | -2.15 | -0.27 | 0.02 | 0.64 | 0.26 | 2.24 | 0.00 | -0.35 | 0.00 |
| 154 | benzidine | Prediction | 3.77 | 1.56 | 0.03 | 2.08 | 0.19 | 11.44 | 2.17 | 0.35 | 5.90 |
| 155 | 3,3'-diaminobenzidine | Prediction | -2.70 | -1.75 | 0.06 | 2.12 | 0.18 | 8.15 | 2.18 | 2.37 | 6.19 |
| 156 | 3,3'-dichlorobenzidine | Prediction | -2.00 | -1.12 | 0.13 | 1.52 | 0.18 | 5.07 | 3.15 | 1.33 | 6.19 |
| 157 | 4,4'-diamino-3,3'-dimethoxy-1,1'-biphenyl | Training | -0.05 | -0.98 | 0.08 | 1.68 | 0.17 | 6.28 | 2.54 | 0.64 | 6.35 |
| 158 | 3,3'-dimethyl-4,4'-bianiline | Training | -1.24 | -1.22 | 0.06 | 1.72 | 0.18 | 7.23 | 2.18 | 0.98 | 6.19 |
| 159 | 4-aminophenyl disulfide | Training | -0.55 | -0.67 | 0.04 | 2.00 | 0.18 | 2.66 | 2.92 | 0.93 | 5.27 |

. Table 4.9.  Continued

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|----|------|--------|----------------------------------|--------------------------------|------------------------|------|----------|------|--------|--------|--------|
| 160 | (4-aminophenyl)ether' | Prediction | -0.92 | -0.41 | 0.03 | 2.00 | 0.18 | 9.41 | 2.46 | -0.04 | 5.67 |
| 161 | 4-aminophenyl Sulfide | Training | 1.02 | -0.22 | 0.02 | 2.00 | 0.18 | 5.92 | 2.43 | 0.26 | 5.67 |
| 162 | 3-OSO3K-AAB | Training | 0.23 | -1.41 | 0.05 | 1.25 | 0.16 | 7.06 | 4.27 | -0.79 | 6.26 |
| 163 | 4'-OSO3K-AAB | Training | -0.75 | -0.62 | 0.06 | 1.42 | 0.16 | 21.87 | 4.32 | 1.40 | 6.30 |
| 164 | 3-methoxy-4'-N,N-diethyl-AAB | Training | -0.99 | -2.08 | 0.03 | 1.54 | 0.15 | 16.01 | 3.87 | 0.26 | 6.25 |
| 165 | 3-OSO3K-MAB | Training | -0.46 | -0.89 | 0.05 | 1.30 | 0.16 | 6.07 | 4.61 | -0.49 | 6.31 |
| 166 | 4'-OSO3K-MAB | Training | -0.89 | -0.66 | 0.08 | 1.38 | 0.16 | 23.36 | 4.59 | -0.71 | 6.32 |
| 167 | 3'-acetoxymethyl-DAB | Training | -1.96 | -2.51 | 0.05 | 1.44 | 0.14 | 6.26 | 4.02 | -0.26 | 6.20 |
| 168 | R1=OCH2CH2OH, R2=NEt2 | Training | -0.38 | -1.80 | 0.14 | 1.41 | 0.09 | 7.89 | 4.15 | -1.09 | 7.83 |
| 169 | R1=OBu, R2=NEt2 | Training | -0.83 | -0.90 | 0.07 | 1.52 | 0.09 | 7.09 | 4.01 | -0.85 | 7.84 |
| 170 | R1=OPr, R2=NEt2 | Training | -0.29 | -0.40 | 0.05 | 1.50 | 0.09 | 7.46 | 4.16 | 0.31 | 7.83 |
| 171 | R1=OMe, R2=NEt2 | Training | -2.15 | -1.09 | 0.04 | 1.45 | 0.09 | 10.97 | 4.05 | 0.18 | 7.81 |
| 172 | R1=OPr, R2=H | Training | -1.43 | -0.64 | 0.07 | 1.33 | 0.10 | 5.34 | 3.92 | 0.08 | 7.78 |
| 173 | R1=OEt, R2=NEt2 | Training | -1.47 | -0.14 | 0.08 | 1.53 | 0.09 | 9.63 | 4.19 | 0.25 | 7.83 |
| 174 | R1=OBu, R2=H | Prediction | -1.01 | -0.80 | 0.09 | 1.36 | 0.09 | 4.14 | 3.80 | -0.28 | 7.79 |

Table 4.9. Continued

| ID | Name | Status | Exp. mutagenicity Log (rev/nmol) | Pred. mutagenicity by model M2 | Hat value ($h^*$=0.15) | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 175 | R1=OCH2CH2OH, R2= H | Training | 0.31 | -0.29 | 0.05 | 1.32 | 0.10 | 6.08 | 3.60 | -0.16 | 7.78 |
| 176 | R1=OMe, R2=H | Training | -1.03 | 0.60 | 0.07 | 1.41 | 0.10 | 6.39 | 3.84 | 0.67 | 7.76 |
| 177 | R1=OEt, R2=H | Training | -1.14 | -0.71 | 0.05 | 1.37 | 0.10 | 5.89 | 3.95 | -1.05 | 7.78 |

**4.3. Prediction of Mutagenicity Activity of External Set Chemicals with No Mutagenicity Data**

The predicted ability of model M2 was tested on external set chemicals with no mutagenic activity data (Insubria graph, Figure 4.4). In this respect, 33 water-soluble textile dyes containing 27 acid and six direct dyes which were provided by Clariant, Huntsman and Dystar were used as an external dataset. The descriptors used in model M2 were calculated for each chemical in the external dataset. Regarding the Insubria graph (Figure 4.4.) the structural coverage of model M2 is 70%.

The chemicals which fall outside the applicability domain mostly belong to azo dyes. These are Direct Blue 78, Acid Red 154, Acid Red 134, Acid Red 138, Acid Red 145, Acid Orange 156, Acid Red 111, Direct Yellow 86, Direct Yellow 106 and one of them is anthraquinone dye (Acid Blue 127). There were many different chemical structures. Because of this heterogenicity of the external set, the 10 compounds were outside the AD. However, 23 compounds belonging both anthraquinone and azo dyes were well predicted by the proposed model (Eq.4.1). The predicted mutagenicity values from model M2 and the values of descriptors for the external set chemicals with no mutagenic activity data are given in Table 4.10.

When the individual descriptor range criterion is considered (Table 4.8.), the chemicals that were out of the descriptor ranges were removed from the list of chemicals given in Table 4.10. The new list of chemicals which fell into the AD of model M2 is given in Table 4.11. The order of the most mutagenic six dyes which fell into the AD of the selected model M2 is as follows:

Acid Blue 62 > Acid Blue 40 > Acid Blue 45> Acid Blue 80 > Acid Blue 230 > Acid Blue 344 These dyes are mainly used for cotton, fiber dyeing and leather shading.

The structures of these dyes are given in Figure 4.5. All these dyes belong to anthraquinone group and they have a substituent in the para position. The most mutagenic dye, Acid Blue 62, has a cyclohexyl group as a substituent. The structure of least mutagenic dye, Direct Orange 34, from the azo group is given in Figure 4.6. It is a synthetic dye that is mainly used for cotton, silk, wool and their blended fabric dyeing and printing, also can be used for leather and paper shading.

Figure 4.4. Insubria graph of model M2 training set in yellow and external set in red.

Table 4.10. The compounds that are in the AD of model M2, their predicted mutagenicity values.

| Chemical Group | Name | CAS Number | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 | Pred. by Eq.(4.1) |
|---|---|---|---|---|---|---|---|---|---|
| Anthraquinone | Acid Blue 40 | 6424-85-7 | 1.184 | 0.085 | 1.710 | 3.498 | 2.380 | 7.597 | -1.212 |
| Anthraquinone | Acid Blue 45 | 2861-02-1 | 1.211 | 0.096 | 5.940 | 2.736 | 2.614 | 7.520 | -1.966 |
| Anthraquinone | Acid Blue 62 | 4368-56-3 | 1.408 | 0.090 | 1.260 | 2.921 | 2.684 | 7.339 | -1.187 |
| Anthraquinone | Acid Blue 80 | 4474-24-2 | 2.170 | 0.074 | 7.000 | 3.934 | -0.158 | 7.866 | -1.659 |
| **\*Azo Dye** | **Acid Blue 113** | 3351-05-1 | **2.310** | 0.073 | 16.050 | 4.43 | **-1.382** | **8.258** | -2.586 |
| Anthraquinone | Acid Blue 230 | 12219-37-3 | 1.172 | 0.086 | 2.470 | 3.756 | 0.520 | 7.571 | -1.962 |
| Anthraquinone | Acid Blue 344 | 85153-93-1 | 1.569 | 0.081 | 4.510 | 3.936 | 0.554 | 7.784 | -1.658 |
| Anthraquinone | Acid Green 25 | 4403-90-1 | 1.783 | 0.075 | 2.060 | 4.49 | 0.960 | 7.765 | -0.662 |
| **Azo Dye** | **Acid Orange 67** | 12220-06-3 | 1.630 | 0.100 | 19.300 | 4.098 | **-1.123** | 7.272 | -3.441 |
| Azo Dye | Acid Orange 127 | 12269-96-4 | 1.819 | 0.094 | 9.750 | 4.435 | -0.223 | 7.612 | -1.837 |
| Azo Dye | Acid Red 42 | 6245-60-9 | 1.449 | 0.096 | 7.070 | 3.913 | -0.717 | 7.729 | -2.574 |
| Azo Dye | Acid Red 57 | 12217-34-4 | 1.366 | 0.094 | 8.870 | 3.802 | -0.454 | 7.730 | -2.768 |
| Azo Dye | Acid Red 151 | 6406-56-0 | 1.426 | 0.101 | 17.930 | 4.367 | 0.630 | 7.058 | -2.474 |
| Azo Dye | Acid Red 249 | 6416-66-6 | 1.425 | 0.086 | 3.040 | 4.403 | -0.431 | 8.157 | -1.844 |
| Azo Dye | Acid Red 266 | 57741-47-6 | 0.760 | 0.106 | 4.650 | 3.423 | 2.776 | 7.573 | -1.711 |
| Azo Dye | Acid Red 337 | 67786-14-6 | 0.919 | 0.106 | 5.740 | 3.041 | 2.703 | 7.552 | -1.959 |
| **Azo Dye** | **Acid Yellow 17** | 6359-98-4 | 0.936 | 0.127 | 15.840 | **5.519** | -0.806 | 6.869 | -2.538 |
| **Azo Dye** | **Acid Yellow 49** | 12239-15-5 | 0.979 | 0.136 | 9.370 | **5.479** | -1.379 | 6.995 | -2.242 |
| Azo Dye | Acid Yellow 61 | 12217-38-8 | 1.199 | 0.102 | 6.870 | 4.235 | -1.078 | 7.177 | -2.635 |
| Azo Dye | Acid Yellow 199 | 70865-20-2 | 1.163 | 0.121 | 8.980 | 3.718 | 1.434 | 6.975 | -2.005 |
| **Azo Dye** | **Direct Blue 71** | 4399-55-7 | 2.201 | 0.056 | 10.740 | **4.938** | **-1.683** | **8.857** | -2.119 |
| **Azo Dye** | **Direct Green 26** | 6388-26-7 | 1.912 | 0.056 | 4.720 | 4.267 | -0.898 | **8.849** | -1.986 |
| Azo Dye | Direct Orange 34 | 12222-37-6 | 1.126 | 0.160 | 16.520 | 4.243 | -0.200 | 6.192 | -2.907 |

**\***Chemicals that are out of the descriptor ranges are written in bold.

Table 4.11. The compounds that are in the AD of model M2 remaining after the application of descriptor range criteria.

| Chemical Group | Compound Name | CIC2 | Chi_D/Dt | L/Bw | TDB09p | Mor28s | piPC08 | Pred. by Eq.(4.1) |
|---|---|---|---|---|---|---|---|---|
| **Anthraquinone*** | **Acid Blue 40** | 1.184 | 2.773 | 1.710 | 3.498 | 2.380 | 7.597 | **1.706** |
| **Anthraquinone** | **Acid Blue 45** | 1.211 | 2.830 | 5.940 | 2.736 | 2.614 | 7.52 | **0.951** |
| **Anthraquinone** | **Acid Blue 62** | 1.408 | 3.077 | 1.260 | 2.921 | 2.684 | 7.339 | **2.060** |
| **Anthraquinone** | **Acid Blue 80** | 2.170 | 2.521 | 7.000 | 3.934 | -0.158 | 7.866 | **0.918** |
| **Anthraquinone** | **Acid Blue 230** | 1.172 | 2.678 | 2.470 | 3.756 | 0.520 | 7.571 | **0.814** |
| **Anthraquinone** | **Acid Blue 344** | 1.569 | 2.350 | 4.510 | 3.936 | 0.554 | 7.784 | **0.754** |
| Azo Dye | Acid Red 151 | 1.426 | 1.885 | 17.930 | 4.367 | 0.630 | 7.058 | -0.678 |
| Azo Dye | Acid Red 266 | 0.760 | 2.154 | 4.650 | 3.423 | 2.776 | 7.573 | 0.484 |
| Azo Dye | Acid Red 337 | 0.919 | 2.210 | 5.740 | 3.041 | 2.703 | 7.552 | 0.278 |
| Azo Dye | Acid Yellow 199 | 1.163 | 1.617 | 8.980 | 3.718 | 1.434 | 6.975 | -0.449 |
| Azo Dye** | Direct Orange 34 | 1.126 | 1.649 | 16.520 | 4.243 | -0.200 | 6.192 | -1.425 |

*The most mutagenic 6 dyes are written in bold; ** The least mutagenic dye.

Figure 4.5. Conformers of (a) Acid Blue 230, (b) Acid Blue 62, (c) Acid Blue 40, (d) Acid Blue 45, (e) Acid Blue 80, (f) Acid Blue 344 drawn in SPARTAN 10 software.

Figure 4.6. A conformer of Direct Orange 34 drawn in SPARTAN 10 software.

## 4.4. Comparison of the QSAR models from the present study with the previously published models

To compare the results of the current study with those of published studies in which mutagenic activity models were generated is of interest, since the comparison points out the strengths and weaknesses of the present study (Table 4.12). Even though an exact comparison is not possible, as each author use different software and a unique dataset with a different number of compounds. Moreover, the quality of prediction depends on various parameters.

Garg et al. (2002) generated several QSAR models on the mutagenicity of aminoazobenzene dyes and related structures and primarily used the best multilinear regression (BMLR) method implemented in CODESSA (CODESSA $^{TM}$, 1992). Later, Bhat et al. (2005) developed several QSAR models on the mutagenicity of aromatic, heteroaromatic amines and related compounds primarily employing again BMLR method implemented in CODESSA (CODESSA $^{TM}$, 1992). They stated that the rather large standard deviation for the test sets suggests the need for improvement in the model. However, they didn't give detail for the standard deviation of test set.

Although $R^2$, $Q^2$, $S$ and $F$ values are provided, these previous works lack the necessary/up-to-date external validation parameters (Table 4.12) However, the fit parameters are not enough to assume that the model is valid and robust. There is no information on the model's applicability domain. Division methods were not given in the previous study. Additionally, training and test set divisions are not clear. All these features reveal that the proposed model in the present study is superior that it passed all of the up-to-date internal and external validation criteria.

Additionally, in attempt to predict the mutagenicity of external chemicals having the most mutagenic activity (Table 4.11) using VEGA software program we uploaded smile code of each chemical to VEGA software. Unfortunately, the output of the program revealed no result for structures with $SO_3H$ (Na) substituents. Therefore, the mutagenic activity values of the most mutagenic chemicals screened from our model could not be verified using VEGA software. Regarding that the VEGA software could not calculate the mutagenic activity of chemicals with sulphonate group, our model is superiror to VEGA software. On the other hand, the least mutagenic dye screened from our QSAR model, Direct Orange 34, was reported as non-mutagenic by VEGA software. This chemical do not contain -$SO_3Na$ group. Therefore, its mutagenicity can be calculated by VEGA. This consistency supports the validation of our model as well. The output of VEGA software for Direct Orange 34 is given in Appendix E.

Table 4.12. Comparison of the statistical parameters of generated models to those of the previously published models.

| Chemical groups | Method | $N^*$ | Number of descriptors | $R^2$ | $Q^2_{LOO}$ | S/F | $R^2_{Test}$ | $RMSE_{Test}$ | Reference |
|---|---|---|---|---|---|---|---|---|---|
| Aminoazodyes | BMLR | 43 | 5 | 0.62 | 0.68 | 0.13/40.40 | N/A | N/A | **Garg et al., 2002** |
| Aromatic/heteroaromatic amine derivatives | BMLR | 181 | 6 | 0.67 | 0.64 | 0.79/58.58 | N/A | N/A | **Bhat et al., 2005** |
| Aromatic/heteroaromatic amine derivatives | MLR | 177 | 6 | 0.61 | 0.56 | 1.04/33.71 | 0.740 | 0.650 | **Present study** |

* the number of chemicals in the datasets.

# 5. CONCLUSION

Aminoazobenzene derivatives and their related structures are really important industrial colorants. However, some of these dyes are shown to be mutagenic.

In the present study, in an attempt to determine the mutagenic activity of azo dyes and amine derivatives numerous QSAR models were generated for different training and test set divisions. All models were validated both internally and externally using recently reported validation parameters in the literature. Their applicability domains were defined by Williams plot and the range of descriptors . The best model from each division was selected via the Multi-Criteria Decision Making (MCDM) approach as implemented in QSARINS. We selected models with the best MCDM score, ensuring the statistical thresholds for fitting, internal and external validation, and with the least possible number of descriptors. The predictive performance of the 5 final models with 6 descriptors from DRAGON software was compared with an external set of 33 compounds mainly comprised of textile dyes. Of the five models, the one with the highest structural coverage (70%) for external set chemicals with no mutagenicity data was proposed as the best model. Descriptors appearing in the proposed model were based on 2D and 3D geometries of the molecule. Mutagenicity is found to be related directly to the abundance of rings in a molecule, polarizability and electric state of a molecule. It is also noteworthy that chemicals with more negatively charged/ionizable group, lipophilic group and carbon-sulfur atom and with large 3D topological distance tended to be more mutagenic. On the other hand, as the values of 3D geometrical descriptor and molecular multiple path count of order 8 increase mutagenic activity decreases.

The generated QSAR model was found to be superior to the previously published literature models regarding the up-to-date validation criteria and OECD principles.

Based on the predicted values, the most and least mutagenic chemicals were screened from the best model generated in the present study. The most mutagenic dyes contain anthraquinone as a common group and have a substituent in the para position. The order of the most mutagenic six dyes are found as Acid Blue 62 > Acid Blue 40 > Acid Blue 45> Acid Blue 80 > Acid Blue 230 > Acid Blue 344. These dyes are mainly used for cotton, fiber dyeing and leather shading. The most mutagenic dye, Acid Blue 62, has a cyclohexyl group as a substituent. The least mutagenic dye is found as Direct Orange 34 which is an azo dye. It was also classified as non-mutagenic by VEGA

software. It is a synthetic dye that is mainly used for cotton, silk, wool and their blended fabric dyeing and printing, also can be used for leather and paper shading.

In the present study, the proposed model has the potential to notify the environmental scientists about the uncertain outcome of chemicals in the environment. The data gap in mutagenicity of currently used some textile dyes were filled. By using the proposed model the dyes manufacturer can have a chance to check the mutagenicity of the new chemicals which fell in the AD of our QSAR model.  In this way, safer and environment-friendly textile dyes can be produced. By using this QSAR model, chemical consumption can be prevented and time can be saved. This is an economic gain in the scientific world, only prioritize chemicals (i.e the most mutagenic chemicals) can be tested, but non-mutagenic chemicals can no longer need to be tested in the laboratory.

# REFERENCES

Ames, B. N., McCann, J., Yamasaki, E. ,1975. Methods for detecting carcinogens and mutagens with the Salmonella/mammalian-microsome mutagenicity test. Mutation Research/Environmental Mutagenesis and Related Subjects, 31 (6), 347-363.

Ames, B. N., Durston, W. E., Yamasaki, E., Lee, F. D., 1973. Carcinogens are mutagens: a simple test system combining liver homogenates for activation and bacteria for detection. Proceedings of the National Academy of Sciences, 70 (8), 2281-2285.

Ames, B. N., Lee, F. D., Durston, W. E., 1973. An improved bacterial test system for the detection and classification of mutagens and carcinogens. Proceedings of the National Academy of Sciences, 70(3), 782-786.

Ames, B. N., 1984. Charles S. Mott prize. The detection of environmental mutagens and potential carcinogens. Cancer, 53 (10), 2034-2040.

Andrew, R. L., 2001. Molecular Modeling Principles and Applications. Prentice Hall, London.

Anupam, K., 1993. Mutagenicity testing of textile dyes with Salmonella/microsome assay. Journal of Environmental Biology, 14, 327-333.

Bafana, A., Devi, S. S., Chakrabarti, T., 2011. Azo dyes: past, present and the future. Environmental Reviews, 19(NA), 350-371.

Banat, I. M., Nigam, P., Singh, D., Marchant, R., 1996. Microbial decolorization of textile-dyecontaining effluents: a review. Bioresource Technology, 58 (3), 217-227.

Basak, S. C., Gute, B. D., 1997. Characterization of molecular structures using topological Indices.SAR and QSAR in Environmental Research, Vol. 7, pp.1–21.

Benigni, R., 2003. Quantitative structure-activity relationship (QSAR) Models of Mutagens and Carcinogens. CRC PRESS, New York, U.S.A.

Benigni, R., Bossa, C., 2011. Mechanisms of chemical carcinogenicity and mutagenicity: a review with implications for predictive toxicology. Chemical Reviews, 111 (4), 2507-2536.

Bhat, K. L., Hayik, S., Sztandera, L., Bock, C. W., 2005. Mutagenicity of aromatic and heteroaromatic amines and related compounds: a QSAR investigation. QSAR and Combinatorial Science, 24 (7), 831-843.

Bi, W., Hayes, R. B., Feng, P., Qi, Y., You, X., Zhen, J., Qu, B., Fu, Z., Chen, M., Chien, H. T. C., 1992. Mortality and incidence of bladder cancer in benzidine-exposed workers in China. American Journal of Industrial Medicine, 21 (4), 481-489.

Caballero, J., Fernández, M., 2006. Linear and nonlinear modeling of antifungal activity of some heterocyclic ring derivatives using multiple linear regression and Bayesian-regularized neural networks. Journal of Molecular Modeling, 12, 168-181.

Carneiro, P. A., Umbuzeiro, G. A., Oliveira, D. P., Zanoni, M. V. B., 2010. Assessment of water contamination caused by a mutagenic textile effluent/dyehouse effluent bearing disperse dyes. Journal of Hazardous Materials, 174 (1-3), 694-699.

Carreón, T., Hein, M. J., Viet, S. M., Hanley, K. W., Ruder, A. M., Ward, E. M., 2010. Increased bladder cancer risk among workers exposed to o-toluidine and aniline: a reanalysis. Occupational and Environmental Medicine, 67 (5), 348-350.

Carroll, C. C., Warnakulasuriyarachchi, D., Nokhbeh, M. R., Lambert, I. B., 2002. Salmonella typhimurium mutagenicity tester strains that overexpress oxygen-insensitive nitroreductases nfsA and nfsB. Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis, 501 (1-2), 79-98.

Chen, X., Deng, Q., Lin, S., Du, C., Zhao, S., Hu.,Y., Yang Z., Lyu Y., Han, J., 2017. A new approach for risk assessment of aggregate dermal exposure to banned azo dyes in textiles. Regulatory Toxicology and Pharmacology, 91, 173-178.

Chirico, N., Gramatica, P., 2011. Real external predictivity of QSAR models: how to evaluate it. Comparison of different validation criteria and proposal of using the concordance correlation coefficient. Journal of Chemical Information and Modeling, 51 (9), 2320-2335.

Chirico, N., Gramatica, P., 2012. Real external predictivity of QSAR models. Part 2. New intercomparable thresholds for different validation criteria and the need for scatter plot inspection. Journal of Chemical Information and Modeling, 52 (8), 2044-2058.

Chung, K. T., Cerniglia, C. E., 1992. Mutagenicity of azo dyes: structure-activity relationships. Mutation Research/Reviews in Genetic Toxicology, 277 (3), 201-220.

Chung, K. T., Fulk, G. E., Andrews, A. W., 1978. The mutagenicity of methyl orange and metabolites produced by intestinal anaerobes. Mutation Research/Genetic Toxicology, 58(2-3), 375-379.

Chung, K. T., Chen, S. C., Claxton, L. D., 2006. Review of the Salmonella typhimurium mutagenicity of benzidine, benzidine analogues, and benzidine-based dyes. Mutation Research/Reviews in Mutation Research, 612 (1), 58-76.

Claxton, L. D., Allen, J., Auletta, A., Mortelmans, K., Nestmann, E., Zeiger, E., 1987. Guide for the Salmonella typhimurium/mammalian microsome tests for bacterial mutagenicity. Mutation Research/Genetic Toxicology, 189 (2), 83-91.

Claxton, L. D., Houk, V. S., Hughes, T. J., 1998. Genotoxicity of industrial wastes and effluents. Mutation Research/Reviews in Mutation Research, 410 (3), 237-243.

CODESSA $^{TM}$ V.2.0, 1992. Semichem, Florida, U.S.A.

Consonni, V., Ballabio, D., Todeschini, R., 2009. Comments on the definition of the Q 2 parameter for QSAR validation. Journal of Chemical Information and Modeling, 49 (7), 1669-1678.

Consonni, V., Ballabio, D., Todeschini, R., 2010. Evaluation of model predictive ability by external validation techniques. Journal of Chemometrics, 24 (3-4), 194-201.

Couto, S. R., 2009. Dye removal by immobilised fungi. Biotechnology Advances, 27 (3), 227-235.

Cox, P., Pedersen, T., 2002. Directive 2002/61/EC of the European Parliament and of the Council amending for the nineteenth time Council Directive 76/769/EEC relating to restrictions on the marketing and use of certain dangerous substances and preparation (azocolourants). Official Journal of the European Union L 243. 15-18.

De Aragão Umbuzeiro, G., Freeman, H., Warren, S. H., Kummrow, F., Claxton, L. D., 2005. Mutagenicity evaluation of the commercial product CI Disperse Blue 291 using different protocols of the Salmonella assay. Food and Chemical Toxicology, 43 (1), 49-56.

De Campos Ventura-Camargo, B., Marin-Morales, M. A., 2013. Azo dyes: characterization and toxicity-a review. Textiles and Light Industrial Science and Technology, 2 (2), 85-103.

Devillers, J., Mombelli, E., 2010. Evaluation of the OECD QSAR Application Toolbox and Toxtree for estimating the mutagenicity of chemicals. Part 1. Aromatic amines. SAR and QSAR in Environmental Research, 21 (7-8), 753-769.

Dimitrov, S., Dimitrova, G., Pavlov, T., Dimitrova, N., Patlewicz, G., Niemela, J., Mekenyan, O., 2005. A stepwise approach for defining the applicability domain of SAR and QSAR models. Journal of Chemical Information and Modeling, 45 (4), 839-849.

Egli, R., Peter, A. P., Freeman, H. S., 1991. Colour chemistry: the design and synthesis of organic dyes and pigments. Elseiver, London.

Esancy, J. F., Freeman, H. S., Claxton, L. D., 1990. The effect of alkoxy substituents on the mutagenicity of some aminoazobenzene dyes and their reductive-cleavage products. Mutation Research/Reviews in Genetic Toxicology, 238 (1), 1-22.

Forgacs, E., Cserhati, T., Oros, G., 2004. Removal of synthetic dyes from wastewaters: a review. Environment International, 30 (7), 953-971.

Freeman, H. S., 2013. Aromatic amines: use in azo dye chemistry. Frontiers in bioscience (Landmark Ed), 18, 145-164.

Garg, A., Bhat, K. L., Bock, C. W., 2002. Mutagenicity of aminoazobenzene dyes and related structures: a QSAR/QPAR investigation. Dyes and Pigments, 55 (1), 35-52.

Golbraikh, A., Shen, M., Xiao, Z., Xiao, Y. D., Lee, K. H., Tropsha, A., 2003. Rational selection of training and test sets for the development of validated QSAR models. Journal of Computer-aided Molecular Design, 17 (2-4), 241-253.

Gramatica, P., Chirico, N., Papa, E., Cassani, S., Kovarich, S., 2013. QSARINS: A new software for the development, analysis, and validation of QSAR MLR models. Journal of Computational Chemistry, 34 (24), 2121-2132.

Halder, A. K., Jha, T., 2010. Validated predictive QSAR modeling of N-aryl-oxazolidinone-5-carboxamides for anti-HIV protease activity. Bioorganic and Medicinal Chemistry Letters, 20 (20), 6082-6087.

Hashimoto, Y., Watanabe, H., Degawa, M., 1977. Mutagenicity of methoxyl derivatives of N-hydroxy-4-aminoazobenzene and 4-nitroazobenzene. GANN Japanese Journal of Cancer Research, 68 (3), 373-374.

Holkar, C. R., Jadhav, A. J., Pinjari, D. V., Mahamuni, N. M., Pandit, A. B., 2016. A critical review on textile wastewater treatments: possible approaches. Journal of Environmental Management, 182, 351-366.

https://www.vegahub.eu/community/developers;VEGA Developers: Istituto di Ricerche Farmacologiche Mario Negri IRCCS. Date accessed June 2019.

Hubble, M. A., Beck, K. R., O'Neal, W. G., Sharma, Y. C., 2012. Celluiosic substrates for removal ofpollutants from aqueous systems: A review. 2. Dyes. BioResources, 7 (2), 2592-2687.

IARC Working Group on the Evaluation of Carcinogenic Risks to Humans., 2010. Carbon black, titanium dioxide, and talc. IARC monographs on the evaluation of carcinogenic risks to humans, 93,1.

IBM, SPSS v. 25, IBM Corp, Armonk, NY. https://www.ibm.com/us-en/marketplace?lnk=mp. Date accessed January 2019.

Jin, X. C., Liu, G. Q., Xu, Z. H., Tao, W. Y., 2007. Decolorization of a dye industry effluent by Aspergillus fumigatus XC6. Applied Microbiology and Biotechnology, 74 (1), 239-243.

Kaur, A., Sandhu, R. S., Grover, I. S., 1993. Screening of azo dyes for mutagenicity with Ames/Salmonella assay. Environmental and Molecular Mutagenesis, 22 (3), 188-190.

Khalid, A., Arshad, M., Crowley, D. E., 2009. Biodegradation potential of pure and mixed bacterial cultures for removal of 4-nitroaniline from textile dye wastewater. Water Research, 43 (4), 1110-1116.

Kode srl, Dragon (software for molecular descriptor calculation) version 7.0.10, 2017, https://chm.kode-solutions.net. Date accessed January 2019.

Kojima, M., Degawa, M., Hashimoto, Y., Tada, M., 1991. Different effects of DNA adducts induced by carcinogenic and noncarcinogenic azo dyes on in vitro DNA synthesis. Biochemical and Biophysical Research Communications, 179 (2), 817-823.

Kothari V., Joshi S., 2013. Product development and role of merchandiser. Bangladesh Textile Today http://www.textiletoday.com.bd/magazine/212; (Date accessed December 2013).

Kušić, H., Rasulev, B., Leszczynska, D., Leszczynski, J., Koprivanac, N., 2009. Prediction of rate constants for radical degradation of aromatic pollutants in water matrix: A QSAR study. Chemosphere, 75 (8), 1128-1134.

Lawrence, I., Lin, K., 1989. A concordance correlation coefficient to evaluate reproducibility. Biometrics, 45 (1), 255-268.

Lawrence, I., Lin, K., 1992. Assay validation using the concordance correlation coefficient. Biometrics, 48 (1), 599-604.

Lin, G. H., Solodar, W. E., 1988. Structure-activity relationship studies on the mutagenicity of some azo dyes in the Salmonella/microsome assay. Mutagenesis, 3 (4), 311-315.

Lin, Y. H., Leu, J. Y., 2008. Kinetics of reactive azo-dye decolorization by Pseudomonas luteola in a biological activated carbon process. Biochemical Engineering Journal, 39 (3), 457-467.

Locke, P., Goldberg, A. M., 2006. Chapter 10. Toxicology and the Regulatory Process: Alternatives to Animals in Toxicology, 245-260, CRC PRESS, US.

Luan, F., Cordeiro, M. N. D. S., 2012. Overview of QSAR modelling in rational drug design. Recent trends on QSAR in the pharmaceutical perceptions. 1st edition. Sharjah UAE:Bentham Science Publishers, 194-241.

Luan, F., Xu, X., Liu, H., and Cordeiro, M. N. D. S., 2013. Review of quantitative structure-activity/property relationship studies of dyes: recent advances and perspectives. Coloration Technology, 129 (3), 173-186.

Maas, R., Chaudhari, S., 2005. Adsorption and biological decolourization of azo dye Reactive Red 2 in semicontinuous anaerobic reactors. Process Biochemistry, 40 (2), 699-705.

Mathur, N., Bhatnagar, P., Sharma, P., 2012. Review of the mutagenicity of textile dye products. Universal Journal of Environmental Research and Technology. Volume 2, 1-18.

Meal, P. F., Cocker, J., Wilson, H. K., Gilmour, J. M., 1981. Search for benzidine and its metabolites in urine of workers weighing benzidine-derived dyes. Occupational and Environmental Medicine, 38 (2), 191-193.

Modi, H. A., Rajput, G., Ambasana, C., 2010. Decolorization of water soluble azo dyes by bacterial cultures, isolated from dye house effluent. Bioresource Technology, 101 (16), 6580-6583.

Mombelli, E., Raitano, G., Benfenati, E., 2016. In silico prediction of chemically induced mutagenicity: how to use QSAR models and interpret their results. In  Silico Methods for Predicting Drug Toxicity , 87-105. Humana Press, New York.

Mortelmans, K., Zeiger, E., 2000. The Ames Salmonella/microsome mutagenicity assay. Mutation Research - Fundamental and Molecular Mechanisms of Mutagenesis, 455 (1-2), 29-60.

Myslak, Z. W., Bolt, H. M., Brockmann, W., 1991. Tumors of the urinary bladder in painters: A case-control study. American Journal of Industrial Medicine, 19 (6), 705-713.

Netzeva, T. I., Worth, A. P., Aldenberg, T., Benigni, R., Cronin, M. T., Gramatica, P.,  Myatt, G., 2005. Current status of methods for defining the applicability domain of (quantitative) structure-activity relationships: The Report and Recommendations of ECVAM Workshop 52[1,2]. Alternatives to Laboratory Animals, 33 (2), 155-173.

Nwokonkwo, D.C., 2013. Synthesis of 2-(1,3-Dihydro-3-Oxo-2h-Pyridylpyrr-2-Ylidene)-1, 2-Dihydro- 3h- Pyridylpyrrol- 3- One. IOSR Journal of Applied Chemistry (IOSR-JAC), 4 (6), 74–78.

OECD, 2007. Guidance Document on the Validation of (Quantitative) Structure-Activity Relationship [(Q) SAR] Models. Organisation for Economic Co-operation and Development, Paris, France.

Ogugbue, C. J., Sawidis, T., 2011. Bioremediation and detoxification of synthetic wastewater containing triarylmethane dyes by Aeromonas hydrophila isolated from industrial effluent. Biotechnology Research International. 2011. 11.

Ojha, P. K., Mitra, I., Das, R. N., Roy, K., 2011. Further exploring rm2 metrics for validation of QSPR models. Chemometrics and Intelligent Laboratory Systems, 107 (1), 194-205.

Örücü, E., Tugcu, G., Saçan, M. T., 2014. Molecular structure–adsorption study on current textile dyes. SAR and QSAR in Environmental Research, 25 (12), 983-998.

Pereira, L., Alves, M., 2012. Dyes environmental impact and remediation. In Environmental Protection Strategies for Sustainable Development, 111-162, Springer, Dordrecht.

Pinheiro, H. M., Touraud, E., Thomas, O., 2004. Aromatic amines from azo dye reduction: status review with emphasis on direct UV spectrophotometric detection in textile industry wastewaters. Dyes and Pigments, 61 (2), 121-139.

Prasad, A. A., Rao, K. B., 2010. Physico chemical characterization of textile effluent and screening for dye decolorizing bacteria. Global Journal of Biotechnology and Biochemistry, 5(2), 80-86.

Przystaś, W., Zabłocka-Godlewska, E., Grabińska-Sota, E., 2012. Biological removal of azo and triphenylmethane dyes and toxicity of process by-products. Water, Air, and Soil Pollution, 223 (4), 1581-1592.

ECHA, 2017. Report on the current status of regulatory applicability of non-animal approaches under the REACH, CLP and Biocidal Products regulations.

Robinson, T., Chandran, B., Nigam, P., 2002. Removal of dyes from a synthetic textile dye effluent by biosorption on apple pomace and wheat straw. Water Research, 36 (11), 2824-2830.

Roy, K., Kar, S., Das, R. N., 2015. Statistical methods in QSAR/QSPR. In A primer on QSAR/QSPR modeling, 37-59, Springer, Switzerland.

Saíz-Urra, L., González, M. P., Teijeira, M., 2006. QSAR studies about cytotoxicity of benzophenazines with dual inhibition toward both topoisomerases I and II: 3D-MoRSE descriptors and statistical considerations about variable selection. Bioorganic & Medicinal Chemistry, 14 (21), 7347-7358.

Sanchez, P. S., Sato, M. I., Paschoal, C. M., Alves, M. N., Furlan, E. V., Martins, M. T., 1988. Toxicity assessment of industrial effluents from S. Paulo state, Brazil, using short-term microbial assays. Toxicity Assessment, 3 (1), 55-80.

Saratale, R. G., Saratale, G. D., Chang, J. S., Govindwar, S. P., 2011. Bacterial decolorization and degradation of azo dyes: a review. Journal of the Taiwan Institute of Chemical Engineers, 42 (1), 138-157.

Sarnaik, S., Kanekar, P., 1995. Bioremediation of colour of methyl violet and phenol from a dye-industry waste effluent using Pseudomonas spp. isolated from factory soil. Journal of Applied Bacteriology, 79(4), 459-469.

Schüürmann, G., Ebert, R. U., Chen, J., Wang, B., Kühne, R., 2008. External validation and prediction employing the predictive squared correlation coefficient. Test set activity mean vs training set activity mean. Journal of Chemical Information and Modeling, 48 (11), 2140-2145.

Shahin, M. M., 1985. Mutagenicity evaluation of nitroanilines and nitroaminophenols in *Salmonella typhimurium*. International Journal of Cosmetic Science, 7 (6), 277-289.

Shahin, M. M., 1987. Relationships between structure and mutagenic activity of environmental chemicals. Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis, 181 (2), 243-256.

Shahin, M. M., Bugaut, A., Kalopissis, G., 1980. Structure-activity relationship within a series of m-diaminobenzene derivatives. Mutation Research/Genetic Toxicology, 78 (1), 25-31.

Shahin, M. M., Bugaut, A., Kalopissis, G., 1983. Mutagenicity of nitro-para-phenylenediamine and derivatives in the salmonella-typhimurium mammalian microsome test. In Environmental Mutagenesis. 5 (3), 467-467, John Wiley and Sons Inc, New York. USA.

Shaul, G. M., Holdsworth, T. J., Dempsey, C. R., Dostal, K. A., 1991. Fate of water soluble azo dyes in the activated sludge process. Chemosphere, 22 (1-2), 107-119.

Sheridan, R. P., Feuston, B. P., Maiorov, V. N., Kearsley, S. K., 2004. Similarity to molecules in the training set is a good discriminator for prediction accuracy in QSAR. Journal of Chemical Information and Computer Sciences, 44 (6), 1912-1928.

Shi, L. M., Fang, H., Tong, W., Wu, J., Perkins, R., Blair, R. M. Branham, W.S., Dial, S.L.; Moland, C.L. Sheehan, D. M., 2001. QSAR models using a large diverse set of estrogens. Journal of Chemical Information and Computer Sciences, 41(1), 186-195.

Schuur, J. H., Selzer, P., Gasteiger, J., 1996. The coding of the three-dimensional structure of molecules by molecular transforms and its application to structure-spectra correlations and studies of biological activity. Journal of Chemical Information and Computer Sciences, 36, 334-344.

Song, F., Ge, C., Yang, H., Si, H., 2015. QSAR model for the LC50 of nitrobenzenes. Industrial and Engineering Chemistry Research, 52 (10), 3543-3562.

Sorahan, T., 2008. Bladder cancer risks in workers manufacturing chemicals for the rubber industry. Occupational Medicine, 58 (7), 496-501.

Stead, C. V., 1990. Chemistry of azo colorants. In Colorants and Auxiliaries, 1 (7), 146-195, Society of Dyers and Colourists, UK.

Sztandera, L., Garg, A., Hayik, S., Bhat, K. L., Bock, C. W., 2003. Mutagenicity of aminoazo dyes and their reductive-cleavage metabolites: a QSAR/QPAR investigation. Dyes and Pigments, 59 (2), 117-133.

Todeschini, R., Consonni, V., Mauri, A., Pavan, M., 2004. Detecting "bad" regression models: multicriteria fitness functions in regression analysis. Analytica Chimica Acta, 515 (1), 199-208.

Todeschini, R., Consonni, V., 2009. Molecular descriptors for chemoinformatics: volume I: alphabetical listing/volume II: appendices, references Vol. 41. John Wiley and Sons Inc, New York. USA.

Topliss, J. G., Edwards, R. P., 1979. Chance factors in studies of quantitative structure-activity relationships. Journal of Medicinal Chemistry, 22 (10), 1238-1244.

Trieff, N. M., Biagi, G. L., Sadagopa, V. R., Connor, T. H., Cantelli-Forti, G., Guerra, M. C., Legator, M. S., 1989. Aromatic amines and acetamides in Salmonella typhimurium TA98 and TA100: a quantitative structure-activity relation study. Molecular Toxicology, 2 (1), 53-65.

Tugcu, G., Saçan, M. T., Vracko, M., Novic, M., Minovski, N., 2012. QSTR modelling of the acute toxicity of pharmaceuticals to fish. SAR and QSAR in Environmental Research, 23, 297-310.

Veerasamy, R., Rajak, H., Jain, A., Sivadasan, S., Varghese, C. P., Agrawal, R. K., 2011. Validation of QSAR models-strategies and importance. International Journal of Drug Design and Discovery, 3, 511-519.

Wave Function Inc. Irvine, CA, USA, SPARTAN 10 software for Windows, 2010. (https://www.wavefun.com/), Date accessed June 2019.

Weisburger, J. H., 1997. A perspective on the history and significance of carcinogenic and mutagenic N-substituted aryl compounds in human health. Mutation Research, 376 (1-2), 261.

Weisburger, J. H., 2002. Comments on the history and importance of aromatic and heterocyclic amines in public health. Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis, 506, 9-20.

Wolff, A. H., Oehme, F. W., 1974. Carcinogenic chemicals in food as an environmental health issue. Journal of the American Veterinary Medical Association, 164, 623-629.

Woo, Y. T., Lai, D. Y., 2001. Aromatic Amino and Nitro-Amino Compounds and Their Halogenated Der,ivatives. Patty's Toxicology. Wiley. New York.

Worth, A. P., Van Leeuwen, C. J., Hartung, T., 2004. The prospects for using (Q)SARs in a changing political environment--high expectations and a key role for the European Commission's joint research centre. SAR and QSAR in Environmental Research, 15 (5-6), 331-343.

You, X. Y., Chen, J. G., and Hu, Y. N., 1990. Studies on the relation between bladder cancer and benzidine or its derived dyes in Shanghai. Occupational and Environmental Medicine, 47 (8), 544-552.

Zeiger, E., 1987. Carcinogenicity of mutagens: predictive capability of the Salmonella mutagenesis assay for rodent carcinogenicity. Cancer Research, 47 (5), 1287-1296.

Zhou, Z. X., Liu, Y. H., Zhang, X. L., 2014. Predicting carcinogenicity of anilines by quantitative structure-toxicity relationship. In Applied Mechanics and Materials. 665, 559-562. Trans Tech Publications, Switzerland.

# APPENDIX A: CHEMICAL STRUCTURES OF EXTERNAL SET COMPOUNDS



Figure A1. Acid Blue 113



Figure A2.  Acid Orange 67



Figure A3. Acid Orange 127



Figure A4.  Acid Orange 156

Figure A5. Acid Red 42



Figure A6. Acid Red 57



Figure A7. Acid Red 111

Figure A8. Acid Red 134

Figure A9. Acid Red 138

Figure A10. Acid Red 145

Figure A11. Acid Red 151

Figure A12. Acid Red 154



Figure A13. Acid Red 249



Figure A14. Acid Red 266

Figure A15. Acid Red 337



Figure A16. Acid Yellow 61



Figure A17. Acid Yellow 17

Figure A18. Acid Yellow 49



Figure A19. Acid Yellow 199



Figure A20. Direct Blue 71



Figure A21. Direct Blue 78

Figure A22. Direct Green 26



Figure A23. Direct Orange 34



Figure A24. Direct Yellow 86



Figure A25. Direct Yellow 106

Figure A26. Acid Blue 40



Figure A27. Acid Blue 45



Figure A28. Acid Blue 62

Figure A29. Acid Blue 80
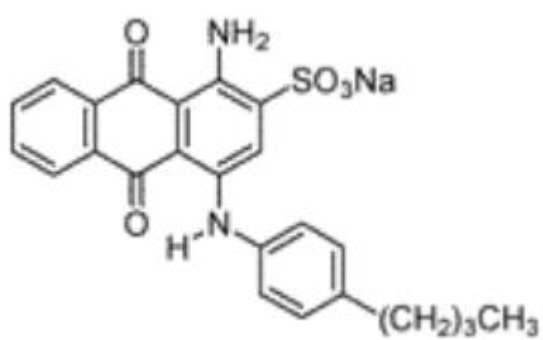


Figure A30. Acid Blue 127



Figure A31. Acid Blue 230
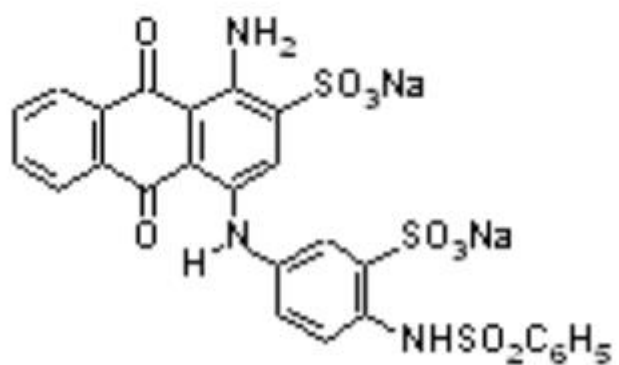
Figure A32. Acid Blue 344

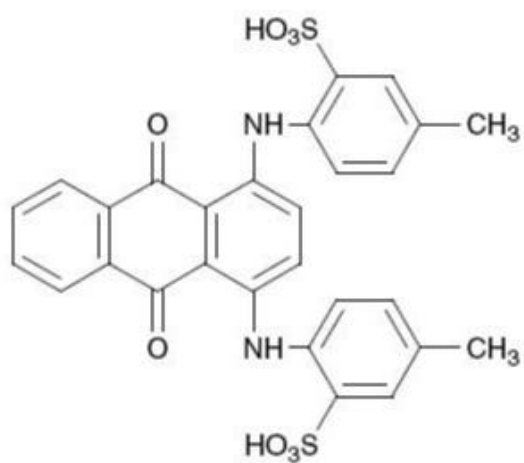

Figure A33. Acid Green 25

# APPENDIX B : PREDICTED VS OBSERVED AND WILLIAMS PLOTS OF THE GENERATED QSAR MODELS
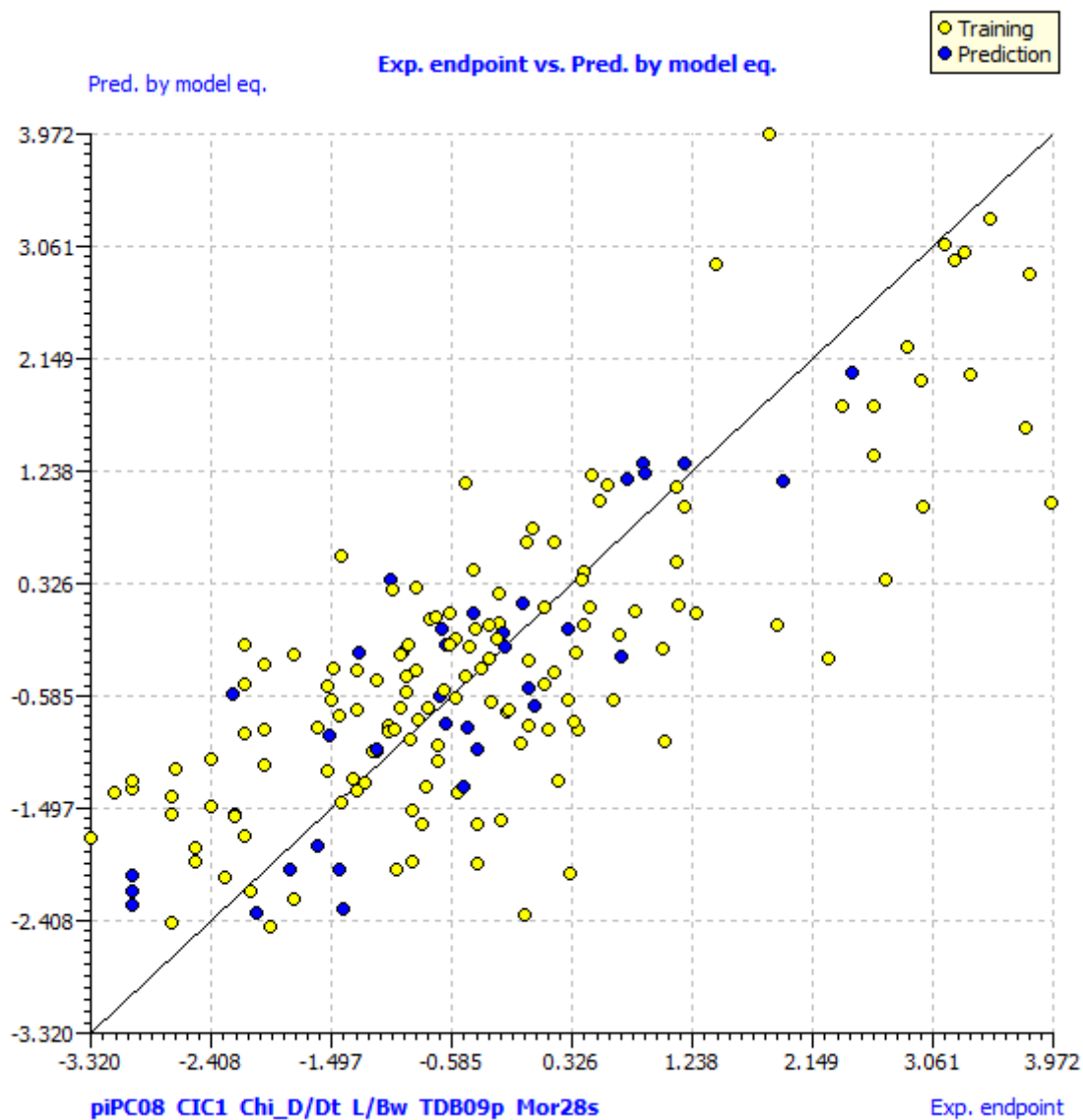


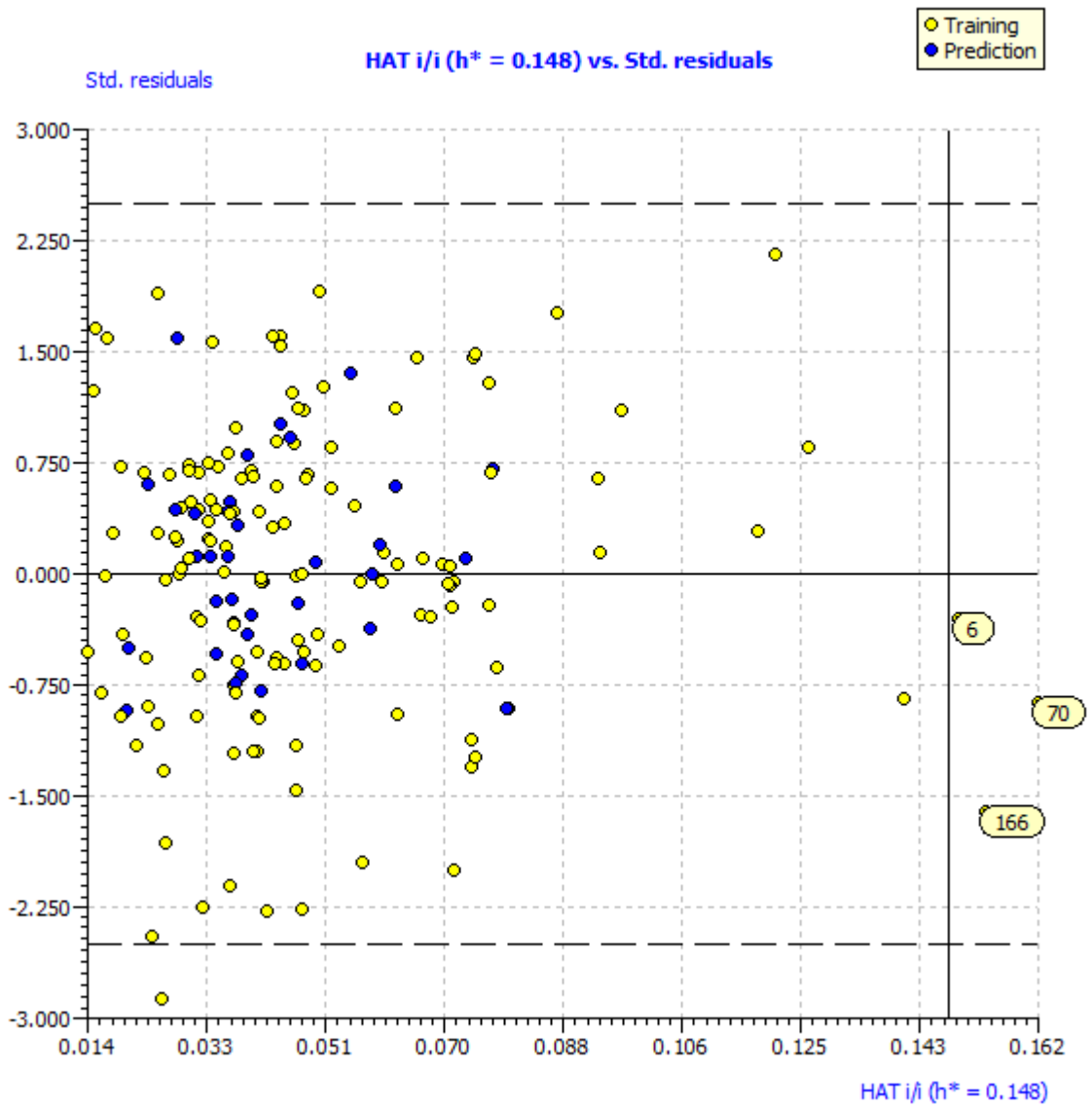Figure B1.  Predicted vs Observed of Model M1

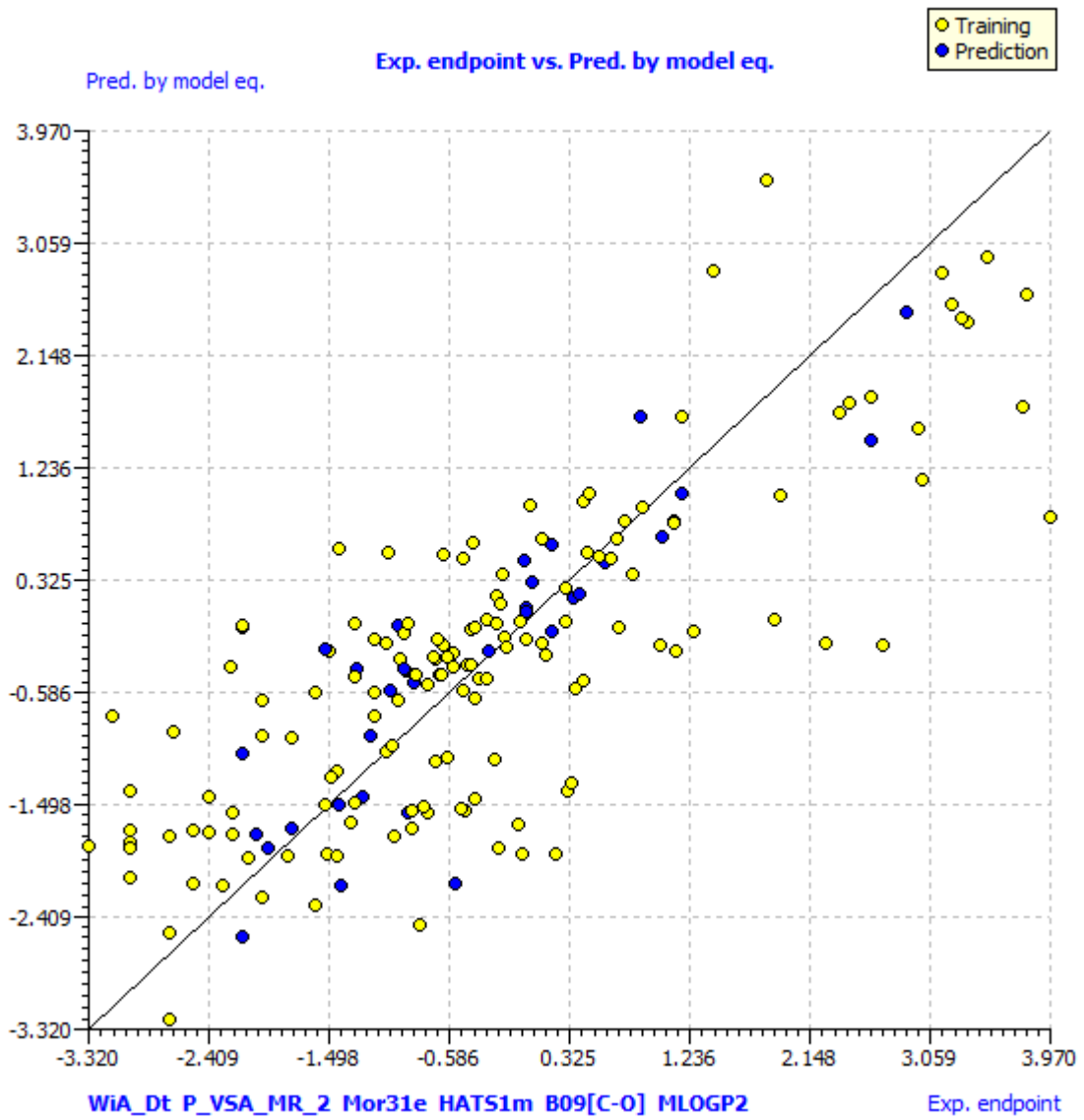Figure B2. Williams plot of Model M1

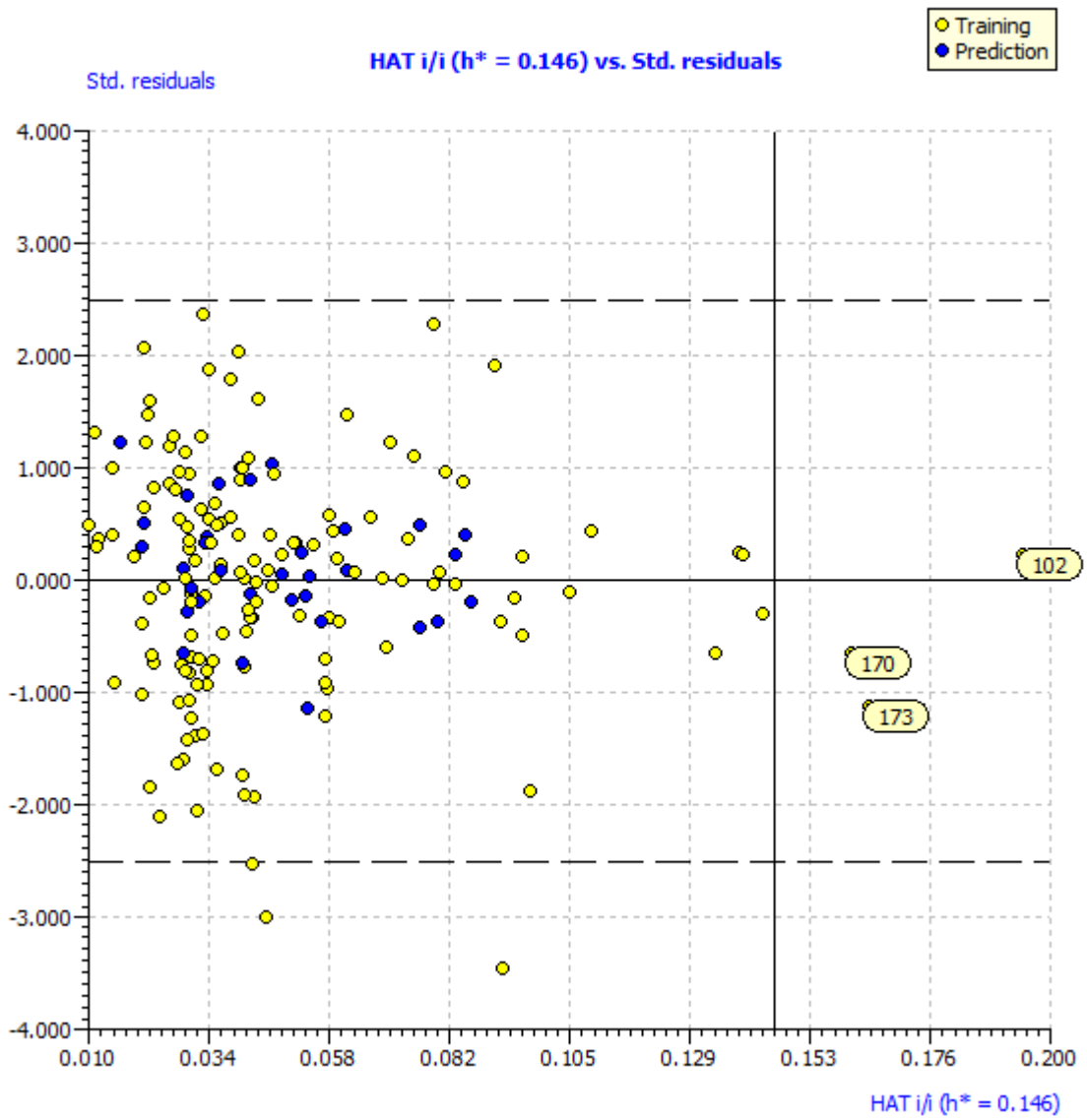Figure B3. Predicted vs Observed of Model M3

Figure B4. Williams plot of Model M3

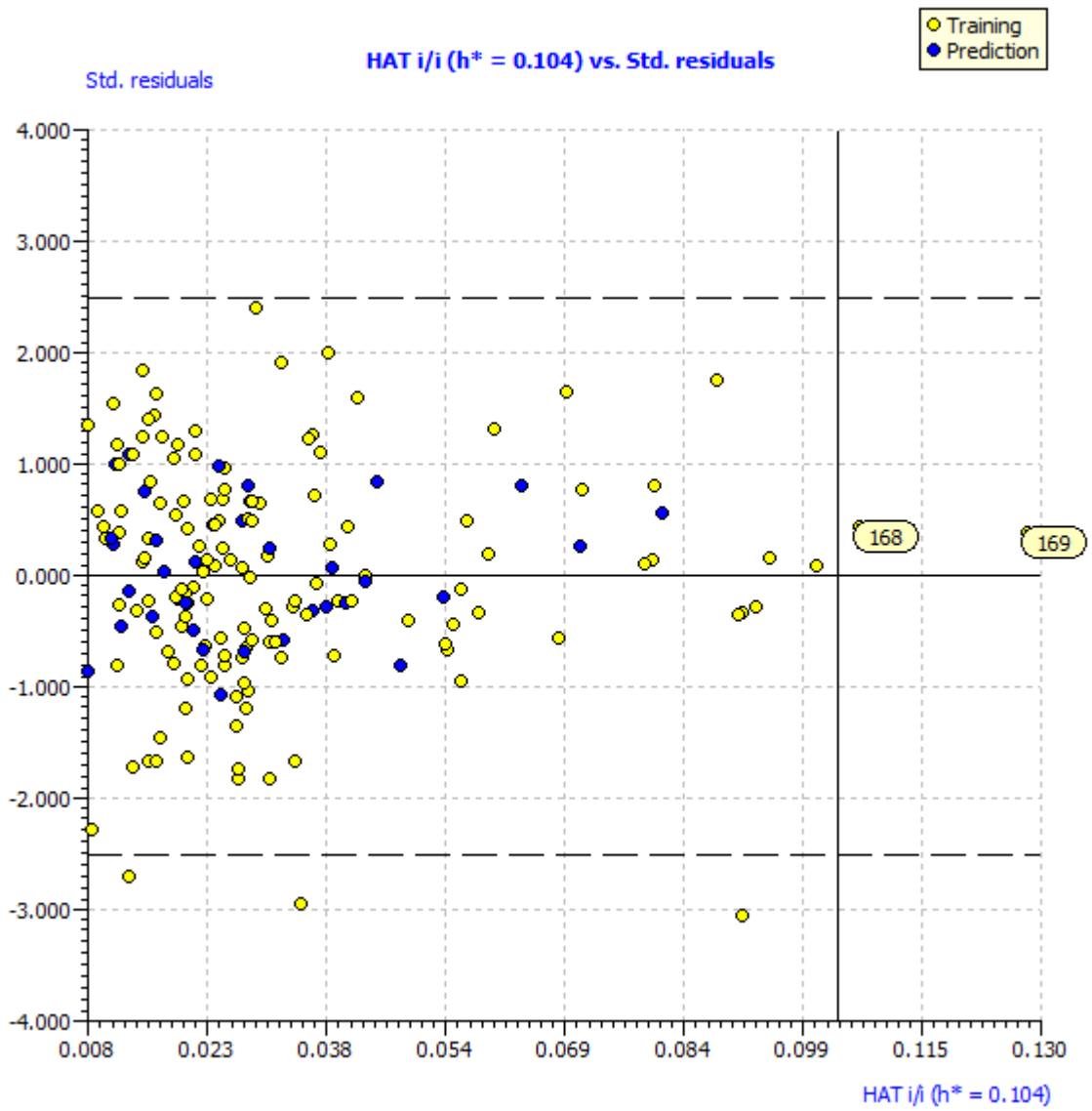Figure B5. Predicted vs Observed of Model M4
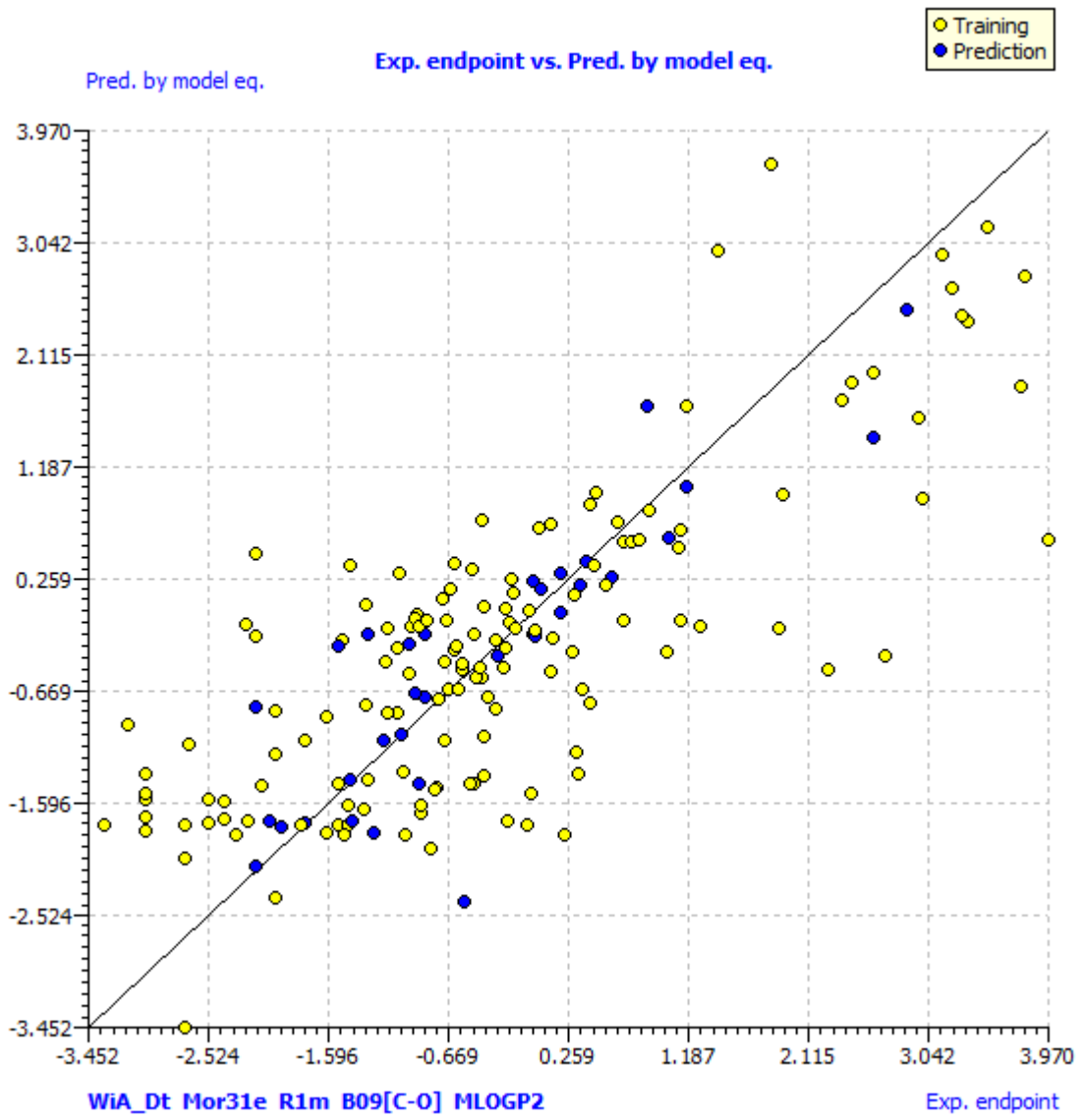
Figure B6. Williams plot of Model M4
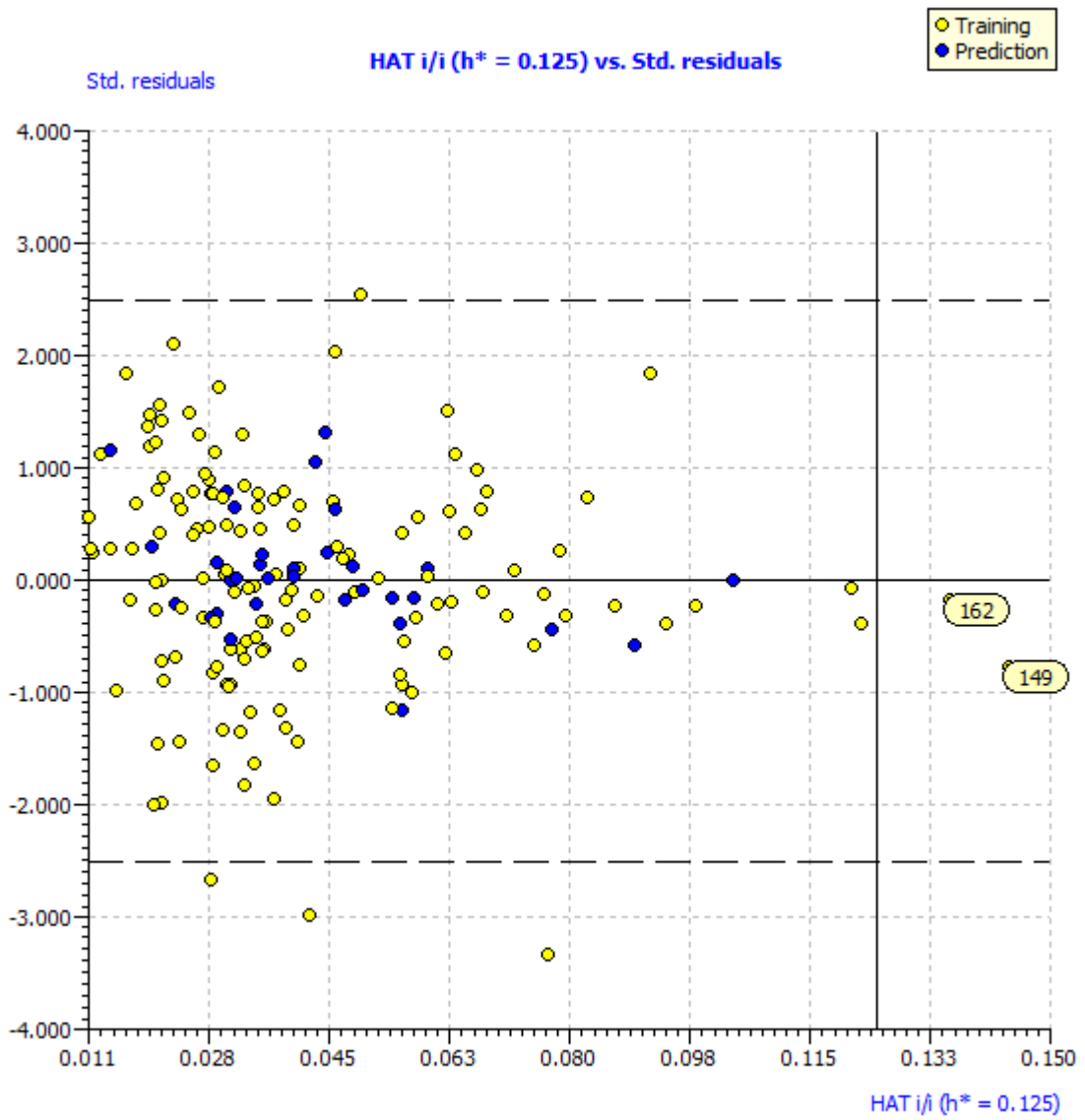
Figure B7. Predicted vs Observed of Model M5

Figure B8. Williams plot of Model M5

# APPENDIX C : THE EQUATIONS OF THE FOUR QSAR MODELS

95% Confidence intervals of descriptor coefficients were given in parentheses.

**M1** **(Eq C.1)**

**Log (mutagenicty)** (TA98+S9 rev/nmol) = -4.9488 (±1.0301) -0.2142 (±0.1526) **piPC08** + 0.8727 (±0.6324)**CIC1** + 1.0999 (±0.2646) **Chi_D/Dt** – 0.0779 (±0.0562) **L/Bw** + 0.7257 (±0.2651) **TDB09p** + 0.5019 (±0.2297) **Mor28s**

**M3** **(Eq C.2)**

**Log (mutagenicty)** (TA98+S9 rev/nmol) = -3.1225 (± 0.6239) + 0.2924 (±0.0948) **WiA_Dt** + 0.0102 (± 0.0065) **P_VSA_MR_2** – 1.8369 (±0.9872) **Mor31e** + 0.7406 (±0.5046) **B09[C-O]** + 0.1378 (±0.0657) **MLOGP2** – 4.6932 (±3.4316)**HATS1m**

**M4** **(Eq C.3)**

**Log (mutagenicty)** (TA98+S9 rev/nmol) = -3.4908 (±0.5318) + 0.3286 (±0.0957)**WiA_Dt** – 0.9550(±0.8994) **Mor31e** + 0.1163 (± 0.0669) **MLOGP2** – 0.0078 (± 0.0059) **P_VSA_MR_2**

**M5** **(Eq C.4)**

 **Log (mutagenicty)** (TA98+S9 rev/nmol) = -2.6199 (±0.6621) + 0.2668 (±0.0804) **WiA_Dt** – 2.5458 (± 0.9550) **Mor31e** – 1.6902 (±0.8417) **R1m+** + 0.8092 (± 0.5323) **B09[C-O]** - 0.1852 (±0.0559) **MLOGP2**

Table C1. Descriptors appeared in the other four QSAR models and their descriptions.

| Abbreviation of Descriptor | Description | Block |
|---|---|---|
| CIC1 | Complementary Information Content index (neighborhood symmetry of 1-order) | Information Indices |
| WiA_Dt | average Wiener-like index from detour matrix | 2D matrix-based descriptors |
| P_VSA_MR_2 | P_VSA-like on Molar Refractivity, bin 2 | P_VSA-like descriptors |
| Mor31e | signal 31 / weighted by Sanderson electronegativity | 3D-MoRSE descriptors |
| B09[C-O] | Presence/absence of C - O at topological distance 9 | 2D Atom Pairs |
| MLOGP2 | squared Moriguchi octanol-water partition coeff. (logP^2) | Molecular properties |
| HATS1m | leverage-weighted autocorrelation of lag 1 / weighted by mass | GETAWAY descriptors |
| R1m+ | R maximal autocorrelation of lag 1 / weighted by mass | GETAWAY descriptors |

# APPENDIX D : OUTPUT of XTERNAL PLUS TOOL FOR QSAR MODELS

Table D1. The results of Xternal Validation Plus tool for model M1

| User Input File Info. | FileName | *MAE*-test M1 .xlsx |
|---|---|---|
| Model biasness test | SystematicErrorResult | Absent |
| | nPE / nNE | 0.7500 |
| | nNE / nPE | 1.3333 |
| | \|MPE / MNE\| | 0.9545 |
| | \|MNE / MPE\| | 1.0477 |
| | AAE - \|AE\| | 0.4456 |
| | $R^2$ (Residuals; serial correlation) | 0.0155 |
| | $R^2$ (Residuals and Yobs values) | 0.1826 |
| | $R^2_{TesT}$ (100% data) | 0.7448 |
| | $R_0{}^2_{TesT}$ (100% data) | 0.7381 |
| | $R_0{}'^2_{TesT}$ (100% data) | 0.7437 |
| Classical Metrics | $Q^2F_1$(100% data) | 0.7531 |
| (for 100% data) | $Q^2F_2$(100% data) | 0.7380 |
| | Scaled Avg.Rm$^2$(100% data) | 0.6488 |
| | Scaled DeltaRm$^2$(100% data) | 0.1754 |
| | *CCC*(100% data) | 0.8567 |
| | $R^2Test$(95% data) | 0.8120 |
| | $R_0{}^2Test$(95% data) | 0.8119 |
| Classical Metric | $R_0{}'^2Test$(95% data) | 0.7550 |
| (after removing | $Q^2F1$(95% data) | 0.8186 |
| 5% data with | $Q^2F2$(95% data) | 0.8115 |
| high residuals) | ScaledAvgRm$^2$(95% data) | 0.7384 |
| | ScaledDeltaRm$^2$(95% data) | 0.1138 |
| | *CCC*(95% data) | 0.8983 |
| | RMSEP(100% data) | 0.6536 |
| Error-based metrics | SD(100% data) | 0.3824 |
| (for 100% data) | SE(100% data) | 0.0646 |
| | *MAE*(100% data) | 0.5340 |
| | RMSEP(95% data) | 0.5561 |
| Error-based metric | SD(95% data) | 0.2963 |
| (after removing 5% data | SE(95% data) | 0.0516 |
| with high residuals) | *MAE*(95% data) | 0.4734 |
| | *MAE*+3*SD(95% data) | 1.3622 |
| BASIC DATA STRUCTURE INFORMATION | | |
| | NCompTest | 35 |
| Number of test set compounds, | Train range | 7.2900 |
| Range and Mean (train and test) | TrainYMean | -0.2900 |
| | Test range | 5.4600 |
| | TestYMean | -0.6063 |
| Distribution of observed | %Y(+/-0.5)TestMean | 31.4286 |
| response values of Test set | %Y(+/-1.0)TestMean | 62.8571 |
| around Test mean(in %) | %Y(+/-1.5)TestMean | 80.0000 |
| | %Y(+/-2.0)TestMean | 85.7143 |
| Distribution of observed | %Y(+/-0.5)TrainMean | 37.1429 |
| response values of Test set | %Y(+/-1.0)TrainMean | 54.2857 |
| around Train mean (in %) | %Y(+/-1.5)TrainMean | 77.1429 |
| | %Y(+/-2.0)TrainMean | 85.7143 |
| | %NComp>(0.1*TR) | 28.5714 |

Table D1.  Continued.

| Distribution of prediction | %NComp>(0.15*TR) | 5.7143 |
|---|---|---|
| errors (in %) | %NComp>(0.2*TR) | 2.8571 |
| | %NComp>(0.25*TR) | 0.0000 |
| | (0.1*TrainingSetRange) | 0.7290 |
| Threshold values utilized | (0.15*TrainingSetRange) | 1.0935 |
| to judge the model predictions | (0.2*TrainingSetRange) | 1.4580 |
| | (0.25*TrainingSetRange) | 1.8225 |
| RESULT (*MAE*-based criteria applied on 95% data) | Prediction Quality | GOOD |

Table D2. The results of Xternal Validation Plus tool for model M3

| User Input File Info. | FileName | *MAE*-test M3 .xlsx |
|---|---|---|
| Model biasness test | SystematicErrorResult | Absent |
| | nPE / nNE | 0.7368 |
| | nNE / nPE | 1.3571 |
| | \|MPE / MNE\| | 0.9801 |
| | \|MNE / MPE\| | 1.0203 |
| | AAE - \|AE\| | 0.3888 |
| | $R^2$ (Residuals; serial correlation) | 0.0139 |
| | $R^2$ (Residuals and Yobs values) | 0.1263 |
| | $R^2_{TesT}$ (100% data) | 0.7671 |
| | $R_0^2{}_{TesT}$ (100% data) | 0.7610 |
| | $R_0'^2{}_{TesT}$ (100% data) | 0.7671 |
| Classical Metrics | $Q^2F_1$(100% data) | 0.7599 |
| (for 100% data) | $Q^2F_2$(100% data) | 0.7588 |
| | Scaled Avg.Rm$^2$(100% data) | 0.6861 |
| | Scaled DeltaRm$^2$(100% data) | 0.1247 |
| | *CCC*(100% data) | 0.8727 |
| | $R^2Test$(95% data) | 0.8437 |
| | $R_0^2Test$(95% data) | 0.8389 |
| Classical Metric | $R_0'^2Test$(95% data) | 0.8230 |
| (after removing | $Q^2F1$(95% data) | 0.8388 |
| 5% data with | $Q^2F2$(95% data) | 0.8386 |
| high residuals) | ScaledAvgRm$^2$(95% data) | 0.7518 |
| | ScaledDeltaRm$^2$(95% data) | 0.1230 |
| | CCC(95% data) | 0.9129 |
| | RMSEP(100% data) | 0.6037 |
| Error-based metrics | SD(100% data) | 0.3928 |
| (for 100% data) | SE(100% data) | 0.0684 |
| | *MAE*(100% data) | 0.4635 |
| | RMSEP(95% data) | 0.5030 |
| Error-based metric | SD(95% data) | 0.3092 |
| (after removing 5% data | SE(95% data) | 0.0555 |
| with high residuals) | *MAE*(95% data) | 0.4006 |
| | *MAE*+3*SD(95% data) | 1.3283 |
| BASIC DATA STRUCTURE INFORMATION | | |
| | NCompTest | 33 |
| Number of test set compounds, | Train range | 7.2900 |
| Range and Mean (train and test) | TrainYMean | -0.3400 |
| | Test range | 5.0300 |
| | TestYMean | -0.4248 |
| Distribution of observed | %Y(+/-0.5)TestMean | 30.3030 |
| response values of Test set | %Y(+/-1.0)TestMean | 63.6364 |
| around Test mean(in %) | %Y(+/-1.5)TestMean | 78.7879 |

Table D2.  Continued.

| | | |
|---|---|---|
| | %Y(+/-2.0)TestMean | 93.9394 |
| Distribution of observed response values of Test set around Train mean (in %) | %Y(+/-0.5)TrainMean | 18.1818 |
| | %Y(+/-1.0)TrainMean | 60.6061 |
| | %Y(+/-1.5)TrainMean | 78.7879 |
| | %Y(+/-2.0)TrainMean | 93.9394 |
| Distribution of prediction errors (in %) | %NComp>(0.1*TR) | 24.2424 |
| | %NComp>(0.15*TR) | 9.0909 |
| | %NComp>(0.2*TR) | 3.0303 |
| | %NComp>(0.25*TR) | 0.0000 |
| Threshold values utilized to judge the model predictions | (0.1*TrainingSetRange) | 0.7290 |
| | (0.15*TrainingSetRange) | 1.0935 |
| | (0.2*TrainingSetRange) | 1.4580 |
| | (0.25*TrainingSetRange) | 1.8225 |
| RESULT (*MAE*-based criteria applied on 95% data) | Prediction Quality | GOOD |

Table D3. The results of Xternal Validation Plus tool for model M4.

| User Input File Info. | FileName | *MAE*-test M4.xlsx |
|---|---|---|
| Model biasness test | SystematicErrorResult | Absent |
| | nPE / nNE | 0.9412 |
| | nNE / nPE | 1.0625 |
| | \|MPE / MNE\| | 0.8561 |
| | \|MNE / MPE\| | 1.1681 |
| | AAE - \|AE\| | 0.4801 |
| | $R^2$ (Residuals; serial correlation) | 0.0083 |
| | $R^2$ (Residuals and Yobs values) | 0.1796 |
| Classical Metrics (for 100% data) | $R^2_{TesT}$ (100% data) | 0.7376 |
| | $R_0^2{}_{TesT}$ (100% data) | 0.7339 |
| | $R_0'^2{}_{TesT}$ (100% data) | 0.7367 |
| | $Q^2F_1$(100% data) | 0.7340 |
| | $Q^2F_2$(100% data) | 0.7327 |
| | Scaled Avg.Rm$^2$(100% data) | 0.6477 |
| | Scaled DeltaRm$^2$(100% data) | 0.1524 |
| | CCC(100% data) | 0.8540 |
| Classical Metric (after removing 5% data with high residuals) | $R^2Test$(95% data) | 0.7779 |
| | $R_0^2Test$(95% data) | 0.7736 |
| | $R_0'^2Test$(95% data) | 0.7694 |
| | $Q_2F1$(95% data) | 0.7733 |
| | $Q_2F2$(95% data) | 0.7719 |
| | ScaledAvgRm2(95% data) | 0.7001 |
| | ScaledDeltaRm2(95% data) | 0.1196 |
| | CCC(95% data) | 0.8793 |
| Error-based metrics (for 100% data) | RMSEP(100% data) | 0.6355 |
| | SD(100% data) | 0.3435 |
| | SE(100% data) | 0.0598 |
| | *MAE*(100% data) | 0.5380 |
| Error-based metric (after removing 5% data with high residuals) | RMSEP(95% data) | 0.5850 |
| | SD(95% data) | 0.3129 |
| | SE(95% data) | 0.0562 |
| | *MAE*(95% data) | 0.4975 |
| | *MAE*+3*SD(95% data) | 1.4362 |
| BASIC DATA STRUCTURE INFORMATION | | |
| Number of test set compounds, Range and Mean (train and test) | NCompTest | 33 |
| | Train range | 7.2900 |
| | TrainYMean | -0.3400 |
| | Test range | 5.0300 |
| | TestYMean | -0.4248 |
| Distribution of observed response values of Test set around Test mean(in %) | %Y(+/-0.5)TestMean | 30.3030 |
| | %Y(+/-1.0)TestMean | 63.6364 |
| | %Y(+/-1.5)TestMean | 78.7879 |
| | %Y(+/-2.0)TestMean | 93.9394 |
| Distribution of observed response values of Test set around Train mean (in %) | %Y(+/-0.5)TrainMean | 18.1818 |
| | %Y(+/-1.0)TrainMean | 60.6061 |
| | %Y(+/-1.5)TrainMean | 78.7879 |
| | %Y(+/-2.0)TrainMean | 93.9394 |
| Distribution of prediction errors (in %) | %NComp>(0.1*TR) | 33.3333 |
| | %NComp>(0.15*TR) | 6.0606 |
| | %NComp>(0.2*TR) | 0.0000 |
| | %NComp>(0.25*TR) | 0.0000 |
| | (0.1*TrainingSetRange) | 0.7290 |

Table D3. Continued.

| | | |
|---|---|---|
| Threshold values utilized | (0.15*TrainingSetRange) | 1.0935 |
| to judge the model predictions | (0.2*TrainingSetRange) | 1.4580 |
| | (0.25*TrainingSetRange) | 1.8225 |
| RESULT (*MAE*-based criteria applied on 95% data) | Prediction Quality | GOOD |

Table D4. The results of Xternal Validation Plus tool for model M5

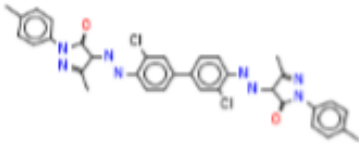| User Input File Info. | FileName | *MAE*-test M5.xlsx |
|---|---|---|
| Model biasness test | SystematicErrorResult | Absent |
| | nPE / nNE | 0.7368 |
| | nNE / nPE | 1.3571 |
| | \|MPE / MNE\| | 1.2017 |
| | \|MNE / MPE\| | 0.8321 |
| | AAE - \|AE\| | 0.4022 |
| | $R^2$ (Residuals; serial correlation) | 0.0120 |
| | $R^2$ (Residuals and Yobs values) | 0.1688 |
| | $R^2Test$(100% data) | 0.7471 |
| | $R_0^2Test$(100% data) | 0.7457 |
| | $R_0'^2Test$(100% data) | 0.7444 |
| Classical Metrics | $Q_2F1$(100% data) | 0.7451 |
| (for 100% data) | $Q_2F2$(100% data) | 0.7439 |
| | Scaled Avg.Rm^2(100% data) | 0.6610 |
| | Scaled DeltaRm^2(100% data) | 0.1236 |
| | *CCC*(100% data) | 0.8609 |
| | $R^2Test$(95% data) | 0.8434 |
| | $R_0^2Test$(95% data) | 0.8426 |
| Classical Metric | $R_0'^2Test$(95% data) | 0.8280 |
| (after removing | $Q_2F1$(95% data) | 0.8417 |
| 5% data with | $Q_2F2$(95% data) | 0.8416 |
| high residuals) | ScaledAvgRm2(95% data) | 0.7491 |
| | ScaledDeltaRm2(95% data) | 0.1248 |
| | *CCC*(95% data) | 0.9122 |
| | RMSEP(100% data) | 0.6220 |
| Error-based metrics | SD(100% data) | 0.4582 |
| (for 100% data) | SE(100% data) | 0.0798 |
| | *MAE*(100% data) | 0.4282 |
| | RMSEP(95% data) | 0.4888 |
| Error-based metric | SD(95% data) | 0.3454 |
| (after removing 5% data | SE(95% data) | 0.0620 |
| with high residuals) | *MAE*(95% data) | 0.3514 |
| | *MAE*+3*SD(95% data) | 1.3876 |
| BASIC DATA STRUCTURE INFORMATION | | |
| | NCompTest | 33 |
| Number of test set compounds, | Train range | 7.2900 |
| Range and Mean (train and test) | TrainYMean | -0.3400 |
| | Test range | 5.0300 |
| | TestYMean | -0.4248 |
| Distribution of observed | %Y(+/-0.5)TestMean | 30.3030 |
| response values of Test set | %Y(+/-1.0)TestMean | 63.6364 |
| around Test mean(in %) | %Y(+/-1.5)TestMean | 78.7879 |
| | %Y(+/-2.0)TestMean | 93.9394 |

Table D4. Continued.

| | | |
|---|---|---|
| Distribution of observed response values of Test set around Train mean (in %) | %Y(+/-0.5)TrainMean | 18.1818 |
| | %Y(+/-1.0)TrainMean | 60.6061 |
| | %Y(+/-1.5)TrainMean | 78.7879 |
| | %Y(+/-2.0)TrainMean | 93.9394 |
| Distribution of prediction errors (in %) | %NComp>(0.1*TR) | 18.1818 |
| | %NComp>(0.15*TR) | 15.1515 |
| | %NComp>(0.2*TR) | 3.0303 |
| | %NComp>(0.25*TR) | 3.0303 |
| Threshold values utilized to judge the model predictions | (0.1*TrainingSetRange) | 0.7290 |
| | (0.15*TrainingSetRange) | 1.0935 |
| | (0.2*TrainingSetRange) | 1.4580 |
| | (0.25*TrainingSetRange) | 1.8225 |
| RESULT (*MAE*-based criteria applied on 95% data) | Prediction Quality | GOOD |

# APPENDIX E: THE OUTPUT OF VEGA SOFTWARE FOR DIRECT ORANGE 34

VEGA

Mutagenicity (Ames test) CONSENSUS model 1.0.3

page 1

## 1. Prediction Summary

**Prediction for compound 1(N.A.)**



Prediction: 🟢

**Prediction is NON-Mutagenic with a consensus score of 0.35, based on 4 models.**

Compound: 1
Compound SMILES:
O=C5N(N=C(C)C5(N=Nc1ccc(cc1Cl)c4ccc(N=NC2C(=O)N(N=C2(C))c3ccc(cc3)C)o(c4)Cl))c6ccc(cc6)C
Used models: 4
Predicted Consensus Mutagen activity: NON-Mutagenic
Mutagenic Score: 0.15
Non-Mutagenic Score: 0.35
Model Caesar assessment: NON-Mutagenic (moderate reliability)
Model ISS assessment: NON-Mutagenic (moderate reliability)
Model SarPy assessment: Mutagenic (moderate reliability)
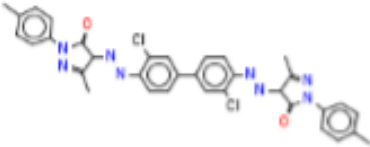Model KNN assessment: NON-Mutagenic (low reliability)
Remarks:
  none

## 1. Prediction Summary

**Prediction for compound 1(N.A.)**

Prediction: 🟢    Reliability: ⭐⭐☆

**Prediction is NON-Mutagenic, but the result shows some critical aspects, which require to be checked:**
**- only moderately similar compounds with known experimental value in the training set have been found**
**- similar molecules found in the training set have experimental values that disagree with the predicted value**

Compound: 1

Compound SMILES:

O=C5N(N=C(C)C5(N=Nc1ccc(cc1Cl)c4ccc(N=NC2C(=O)N(N=C2(C))c3ccc(cc3)C)c(c4)Cl))c6ccc(cc6)C

Experimental value: -

Predicted Mutagen activity: NON-Mutagenic

Structural alerts: -

Reliability: the predicted compound could be out of the Applicability Domain of the model

Remarks: