1LabMol - Laboratory for Molecular Modeling and Drug Design, Faculdade de Farmácia, Universidade Federal de Goiás, Goiânia, Brazil 2Laboratory of Cheminformatics, Centro Universitário de Anápolis (UniEVANGÉLICA), Anápolis, Brazil 3Laboratory for Molecular Modeling, Division of Chemical Biology and Medicinal Chemistry, Eshelman School of Pharmacy, University of North Carolina at Chapel Hill, Chapel Hill, NC, United States 4Department of Chemical Technology, Odessa National Polytechnic University, Odessa, Ukraine

Virtual screening (VS) has emerged in drug discovery as a powerful computational approach to screen large libraries of small molecules for new hits with desired properties that can then be tested experimentally. Similar to other computational approaches, VS intention is not to replace in vitro or in vivo assays, but to speed up the discovery process, to reduce the number of candidates to be tested experimentally, and to rationalize their choice. Moreover, VS has become very popular in pharmaceutical companies and academic organizations due to its time-, cost-, resources-, and labor-saving. Among the VS approaches, quantitative structure-activity relationship (QSAR) analysis is the most powerful method due to its high and fast throughput and good hit rate. As the first preliminary step of a QSAR model development, relevant chemogenomics data are collected from databases and the literature. Then, chemical descriptors are calculated on different levels of representation of molecular structure, ranging from 1D to nD, and then correlated with the biological property using machine learning techniques. Once developed and validated, QSAR models are applied to predict the biological property of novel compounds. Although the experimental testing of computational hits is not an inherent part of QSAR methodology, it is highly desired and should be performed as an ultimate validation of developed models. In this mini-review, we summarize and critically analyze the recent trends of QSAR-based VS in drug discovery and demonstrate successful applications in identifying perspective compounds with desired properties. Moreover, we provide some recommendations about the best practices for QSAR-based VS along with the future perspectives of this approach. Quantitative structure-activity relationship (QSAR) analysis is a ligand-based drug design method developed more than 50 years ago by Hansch and Fujita (1964). Since then and until now, QSAR remains an efficient method for building mathematical models, which attempts to find a statistically significant correlation between the chemical structure and continuous (pIC50, pEC50, Ki, etc.) or categorical/binary (active, inactive, toxic, nontoxic, etc.) biological/toxicological property using regression and classification techniques, respectively (Cherkasov et al., 2014). In the last decades, QSAR has undergone several transformations, ranging from the dimensionality of the molecular descriptors (from 1D to nD) and different methods for finding a correlation between the chemical structures and the biological property. Initially, QSAR modeling was limited to small series of congeneric compounds and simple linear regression methods. Nowadays, QSAR modeling has grown, diversified, and evolved to the modeling and virtual screening (VS) of very large data sets comprising thousands of diverse chemical structures and using a wide variety of machine learning techniques (Cherkasov et al., 2014; Mitchell, 2014; Ekins et al., 2015; Goh et al., 2017). This review is devoted to (i) critical analysis of advantages and disadvantages of QSAR-based VS in drug discovery; (ii) demonstration of several successful QSAR-based discoveries of compounds with desired properties; (iii) description of best practices for the QSAR-based VS; and (iv) discussion of future perspectives of this approach. Best Practices in QSAR Modeling and Validation High-throughput screening (HTS) technologies resulted in the explosion of amount of data suitable for QSAR modeling. As a result, data quality problem became one of the fundamental questions in cheminformatics. As obvious as it seems, various errors in both chemical structure and experimental results are considered as major obstacle to building predictive models (Young et al., 2008; Southan et al., 2009; Williams and Ekins, 2011). Considering these limitations, Fourches et al. (2010; 2015; 2016) developed the guidelines for chemical and biological data curation as a first and mandatory step of the predictive QSAR modeling. Organized into a solid functional process, these guidelines allow the identification, correction, or, if needed, removal of structural and biological errors in large data sets. Data curation procedures include the removal of organometallics, counterions, mixtures, and inorganics, as well as the normalization of specific chemotypes, structural cleaning (e.g., detection of valence violations), standardization of tautomeric forms, and ring aromatization. Additional curation elements include averaging, aggregating, or removal of duplicates to produce a single bioactivity result. Detailed discussion of aforementioned data curation procedures can be found elsewhere (Fourches et al., 2010, 2015, 2016). The Organization for Economic Cooperation and Development (OECD) developed a set of guidelines that the researchers should follow to achieve the regulatory acceptance of QSAR models. According to these principles, QSAR models should be associated with (i) defined end point, (ii) unambiguous algorithm, (iii) defined domain of applicability, (iv) appropriate measures of goodness-of-fit, robustness, and predictivity, and (v) if possible, mechanistic interpretation (OECD, 2004). In our opinion, the additional rule requesting thorough data curation as a mandatory preliminary step to model development should be added there. Continuing Importance of QSAR as Virtual Screening Tool The current pipeline to discover hit compounds in early stages of drug discovery is a data-driven process, which relies on bioactivity data obtained from HTS campaigns (Nantasenamat and Prachayasittikul, 2015). Since the cost of obtaining new hit compounds in HTS platforms is rather high, QSAR modeling has been playing a pivotal role in prioritizing compounds for synthesis and/or biological evaluation. The QSAR models can be used for both hits identification and hit-to-lead optimization. In the latter, a favorable balance between potency, selectivity, and pharmacokinetic and toxicological parameters, which is required to develop a new, safe, and effective drug, could be achieved through several optimization cycles. As no compound need to be synthesized or tested before computational evaluation, QSAR represents a labor-, time-, and cost-effective method to obtain compounds with desired biological properties. Consequently, QSAR is widely practiced in industries, universities, and research centers around the world (Cherkasov et al., 2014). The general scheme of QSAR-based VS approach is shown in Figure 1. Initially, the data sets collected from external sources are curated and integrated to remove or correct inconsistent data. Using these data, QSAR models are developed and validated following OECD guidelines and best practices of modeling. Then, QSAR models are used to identify chemical compounds predicted to be active against selected endpoints from large chemical library (i.e., 105 to 107 chemical structures) is reduced by QSAR models to a smaller number of compounds, which then will be tested experimentally (i.e., 101 to 103 chemical structures) (Kar and Roy, 2013; Tanrikulu et al., 2013). However, it is important to mention that modern VS workflows incorporate additional filtering steps, including: (i) sets of empirical rules [e.g., Lipinski's (Lipinski et al., 1997) rules], (ii) chemical similarity cutoffs, (iii) other QSAR-based filters (e.g., toxicological and pharmacokinetic endpoints), and (iv) chemical feasibility and/or purchasability (Cherkasov et al., 2014). Although the experimental validation of computational hits does not represent part of the QSAR methodology, this should be performed as the final important step. After experimental validation, a multi-parameter optimization (MPO) with QSAR predictions of potency, selectivity, and pharmacokinetic parameters can be conducted. This information will be crucial during hit-to lead and lead optimization design of the compound series, to find the properties balance (potency, selectivity, and PK) related with the effect of different decoration patterns to establish a new series of target compounds for in vivo evaluation. FIGURE 1. QSAR-based virtual screening workflow. QSAR-Based Virtual Screening vs. High-Throughput Screening High-throughput screening (HTS) technologies have rapidly identify large subsets of molecules with desired activity from large screening collections of compounds (105–106 compounds) using automated plate-based experimental assays (Mueller et al., 2012). However, the hit rate of HTS ranges between 0.01% and 0.1% and this highlights the frequently encountered limitation that most of the screened compounds are routinely reported as inactive toward the desired bioactivity (Thorne et al., 2010). Consequently, the drug discovery cost increases according to the number of tested compounds (Butkiewicz et al., 2013). On the other hand, typical hit rates from a validated VS method, including QSAR-based, typically range between 1% and 40%. Thus, VS campaigns are found to have a higher rate of biologically active compounds and at a lower cost than HTS. In this perspective, we show that QSAR-based VS could be used to enrich hit rates of HTS campaigns. For example, Mueller et al. (2010) employed both HTS and QSAR models to search novel positive allosteric modulators for mGlu5, a G-protein coupled receptor involved in disorders like schizophrenia and Parkinson's disease. First, the HTS of approximately 144,000 compounds resulted in a total of 1,356 hits, with a hit rate of 0.94%. Then, this dataset was used to build continuous QSAR models (combining physicochemical descriptors and neural networks), which were subsequently applied to screen a database of approximately 450,000 compounds. Finally, 824 compounds were acquired for biological testing and 232 were confirmed as active (hit rate of 28.2%) (Mueller et al., 2010). In another study, Rodriguez et al. (2010) screened approximately 160,000 compounds to identify 624 antagonists of mGlu5. Further, these data were used to develop QSAR models and, then, applied to screen next 700,000 compounds from ChemDiv database. Among them, 88 of acquired compounds were active, corresponding to a hit rate of 3.6% while the HTS had a hit rate of 0.2% (Mueller et al., 2012). Practical Applications of QSAR-Based Virtual Screening Despite its obvious advantages, QSAR modeling remains underestimated as a VS tool. Unfortunately, QSAR is still seen as a complementary analysis to studies of synthesis and biological evaluation, often introduced in the study without any justification or additional perspective. Despite the small number of VS applications available in the literature, most of them led to the discovery of promising hits and lead candidates. Below, we discuss some successful applications of QSAR-based VS for the discovery of new hits and hit-to-lead optimization. Malaria Malaria is an infectious disease caused by five different species of Plasmodium parasites and transmitted to humans through the bite of infected female mosquitoes of the genus Anopheles. The most lethal species is P. falciparum, which can lead to severe illness and death (Phillips et al., 2017). Malaria is a widespread disease: 91 countries and areas have ongoing transmission. According to World Health Organization (WHO), about 216 million cases and 445,000 deaths from malaria were reported in 2016 (WHO, 2018c). Furthermore, the resistance to antimalarial drugs is a common and growing issue and constitutes a substantial threat for populations in endemic regions (Gorobets et al., 2017; Menard and Dondorp, 2017). In a study reported by Zhang et al. (2013), a data set of 3,133 compounds reported as active or inactive against P. falciparum chloroquine susceptible strain (3D7) was used to develop QSAR models. The models were built using Dragon descriptors (0D, 1D, and 2D), ISIDA-2D fragments descriptors and support vector machines (SVM) method. During QSAR modeling and validation, the data set was randomly divided into modeling and external evaluation set. Additionally, the modeling set was divided multiple times in training and test sets using the Sphere Exclusion algorithm. Then, by using a consensus approach, the QSAR models were applied for VS of the ChemBridge database. After VS, 176 potential antimalarial compounds were identified and submitted to experimental validation along with 42 putative inactive compounds, used as negative controls. Twenty-five compounds presented antimalarial activity in P. falciparum growth inhibition assays and low cytotoxicity in mammalian cells. All 42 compounds predicted as inactives by the models were confirmed experimentally (Zhang et al., 2013). The confirmed experimental hits presented new chemical scaffolds against P. falciparum and could be promising starting points for the development of new optimized antimalarial agents. Schistosomiasis Schistosomiasis is a disease caused by flatworms of the genus Schistosoma that affects 206 million of people worldwide (WHO, 2018d). The current reliance on only one drug, praziquantel, for treatment and control of this disease calls for the urgent discovery of novel anti-schistosomal drugs (Colley et al., 2014). Aiming at discovering new drugs, our group developed binary QSAR models for Schistosoma mansoni thioredoxin glutathione reductase (SmTGR), a validated target for schistosomiasis (Kuntz et al., 2007), to find new structurally dissimilar compounds (Neves et al., 2016). To achieve this goal, we developed a study with the following steps: (i) curation of the largest possible data set of SmTGR inhibitors, (ii) development of rigorously validated and mechanically interpretable models, and (iii) application of generated models for VS of ChemBridge library. Using the QSAR models, we prioritized 29 compounds for further experimental evaluation. As a result, we found that the QSAR models were efficient for discovery of six novel hit compounds active against schistosomula and three hits active against adult worms (hit rate of 20.6%). Among them, 2-[2-(3-methyl-4-nitro-5-isoxazolyl)vinyl]pyridine and 2-(benzylsulfonyl)-1,3-benzothiazole, two compounds representing new chemical scaffolds have activity against schistosomula and adult worms at low micromolar concentrations and therefore represent promising antischistosomal hits for further hit-to-lead optimization (Neves et al., 2016). In another study, we developed continuous QSAR models for a data set of oxadiazoles inhibitors of smTGR (Melo-Filho et al., 2016). Using a combi-QSAR approach, we built a consensus model combining the predictions of individual 2D- and 3D-QSAR models. Then, the model was used for VS of ChemBridge database and the 10 top ranked compounds were further evaluated in vitro against schistosomula and adult worms. Additionally, we applied five highly predictive in-house QSAR models for physico-chemical and biological profiling of these hits. The experimental results showed that 4-nitro-3,5-bis(1-nitro-1H-pyrazol-4-yl)-1H-pyrazole (LabMol-17) and 3-nitro-4-{[(4-nitro-1,2,5-oxadiazol-3-yl)oxy]methyl}-1,2,5-oxadiazole (LabMol-19), two compounds containing new chemical scaffolds (hit rate of 20.6%), were highly active in both life stages of the parasite at low micromolar concentrations (Melo-Filho et al., 2016). Tuberculosis Mycobacterium tuberculosis, the causative agent of tuberculosis (TB), kills about 1.6 million people every year (WHO, 2018e). The current treatment of this disease takes approximately 9 months, which normally leads to noncompliance and, hence, the emergence of multidrug-resistant bacteria (AlMatar et al., 2017). Aiming the discovery of new anti-TB agents, our group used QSAR models to design new series of chalcone (1,3-diaryl-2-propen-1-ones) derivatives. Initially, we retrieved from the literature all chalcone compounds with in vitro inhibition data against M. tuberculosis H37Rv strain. After rigorous data curation, these chalcones were subject to structure-activity relationships (SAR) analysis. Based on SAR rules, bioisosteric replacements were employed to design new chalcone derivatives with optimized anti-TB activity. In parallel, binary QSAR models were generated using several machine learning methods and molecular fingerprints. The fivefold external cross-validation procedure confirmed the high predictive power of the developed models. Using these models, we prioritized series of chalcone derivatives for synthesis and biological evaluation (Gomes et al., 2017). As a result, five 5-nitro-substituted heteroaryl chalcones were found to exhibit MICs at nanomolar concentrations against replicating mycobacteria, as well as low micromolar activity against nonreplicating bacteria. In addition, four of these compounds were more potent than standard drug isoniazid. The series also showed low cytotoxicity against commensal bacteria and mammalian cells. These results suggest that designed heteroaryl chalcones, identified with the help of QSAR models, are promising anti-TB lead candidates (Gomes et al., 2017). Viral Infections Yearly, influenza epidemics can seriously affect all populations in the world. These annual epidemics are estimated to result in about 5 million cases and 650,000 deaths (WHO, 2018b). Influenza virus is mutating constantly, resulting in novel resistant strains, and hence, the development of new anti-influenza drugs active against these new strains is important to prevent epidemics (Laborda et al., 2016). Aiming the discovery of new anti-influenza A drugs, Lian et al. (2015) built binary QSAR models, using SVM and Naïve Bayesian methods, to predict neuraminidase inhibition, a validated protein target for influenza. Then, by using a consensus approach, the QSAR models were applied for VS of ChEMBL database. Among 15,600 compounds screened in an in-house database, 60 compounds were selected to experimental evaluation on neuraminidase activity. Nine inhibitors were identified, five of which were oseltamivir derivatives exhibiting potent neuraminidase inhibition at nanomolar concentrations. Other four active compounds belonged to novel scaffolds, with potent inhibition at low micromolar concentrations (Lian et al., 2015). According to WHO, approximately 35 million people are infected with HIV (WHO, 2018a). The treatment for HIV infections requires a lifelong antiretroviral therapy, targeting different stages of HIV replication cycle. Consequently, because of the emergence of resistance and the lack of tolerability, development of novel anti-HIV drugs is of high demand (Cihlar and Fordyce, 2016; Garbelli et al., 2017). With the purpose of discovering new anti-HIV-1 drugs, Kurczyk et al. (2015) developed a two-step VS approach to prioritize compounds against HIV integrase, an important target to viral replication cycle. The first step was based on binary QSAR models, and the second on privileged fragments. Then, 1.5 million of commercially available compounds were screened, and 13 compounds were selected to experimental evaluation on integrase activity. Nine inhibitors were identified, five of which were oseltamivir derivatives exhibiting potent neuraminidase inhibition at nanomolar concentrations. Other four active compounds belonged to novel scaffolds, with potent inhibition at low micromolar concentrations (Lian et al., 2015). According to WHO, approximately 35 million people are infected with HIV (WHO, 2018a). The treatment for HIV infections requires a lifelong antiretroviral therapy, targeting different stages of HIV replication cycle. Consequently, because of the emergence of resistance and the lack of tolerability, development of novel anti-HIV drugs is of high demand (Cihlar and Fordyce, 2016; Garbelli et al., 2017). With the purpose of discovering new anti-HIV-1 drugs, Kurczyk et al. (2015) developed a two-step VS approach to prioritize compounds against HIV integrase, an important target to viral replication cycle. The first step was based on binary QSAR models, and the second on privileged fragments. Then, 1.5 million of commercially available compounds were screened, and 13 compounds were selected to experimental evaluation. These inhibitors belong to different chemical classes with diverse scaffolds representing potential starting points for the development of new anti-HIV-1 drugs. Central Nervous System (CNS) Disorders Anxiety disorders such as schizophrenia (Nichols and Nichols, 2008; Lacivita et al., 2016). However, the currently marketed drugs targeting 5-HT1A receptor possess severe side effects. To address this, Luo et al. (2014) developed a QSAR-based VS workflow to find new hit compounds targeting 5-HT1A receptor. First, binary QSAR models were generated using Dragon descriptors and several machine learning methods. Then, developed QSAR models were rigorously validated and applied in consensus for VS four commercial chemical databases. Fifteen compounds were selected for experimental testing, and nine of them have proven to be active at low nanomolar concentrations. One of the confirmed hits, [(8a)-6-methyl-9,10-didehydroergolin-8-yl]methanol, showed very high binding affinity (Ki) of 2.3 nM against 5-HT1A receptor. Future Directions and Conclusion To summarize, we would like to emphasize that QSAR modeling represents a time-, labor-, and cost-effective tool to discover hit compounds and lead candidates in the early stages of drug discovery process. Analyzing the examples of QSAR-based VS available in the literature, one can see that many of them led to the identification of promising lead candidates. However, along with success stories, many QSAR projects fail on the model building stage. This is caused by the lack of understanding that QSAR is highly interdisciplinary and application field as well as great experience of the best practices in the field (Tropsha, 2010; Ban et al., 2017). Earlier, we have explained this by the undesirably high population of "button pushers," that is, researchers who conduct modeling without understanding and analyzing the data and modeling process itself (Muratov et al., 2012). This was also explained by the elusive ease of obtaining computational model and making even advanced calculations without understanding of the sense and limitations of the approach (Bajorath, 2012). In addition to this, a lot of even experienced researchers target their efforts to a "vicious statistical cycle," which main goal is to validate models using as many metrics as possible. In this case, the QSAR modeling is restricted to a single simple question: "What is the best metrics or the best statistical method"? Although we recognize that the right choice of statistical approach and especially rigorous external validation are necessary and represent an essential step in any computer-aided drug discovery study, we want to reinforce that QSAR modeling is useful only if it is applied for the solution of a formulated problem and results in development of new compounds with desired properties. As future directions, we would like to point out that the era of big data has just started, and it is still in the chemical/biological data accumulation stage. Therefore, to avoid the situation that the number of assayed compounds available on literature exceeds the modeling capability, the development, and implementation of new machine learning algorithms and data curation methods capable of handling millions of compounds are urgently needed. Finally, the overall success of any QSAR-based VS project depends on the ability of a scientist to think critically and prioritize the most promising hits according to his experience. Moreover, the success rate of collaborative drug discovery projects, where the final selection of computational hits is done by both a modeler and an expert in a given field, is much higher than success rate of the projects driven solely by computational or experimental scientists. Author Contributions All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it for publication. Funding This work was partially funded by Grant No. 1U01CA207160 from NIH and Grant No. 400760/2014-2 from CNPq. CHA is Research Fellow in productivity of CNPq, Conflict of Interest Statement The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest. Acknowledgments The authors would like to thank Brazilian funding agencies, CNPq, CAPES, and FAPEG, for financial support and fellowships. References AlMatar, M., AlMandeal, H., Var, I., Kayar, B., and Köksal, F. (2017). New drugs for the treatment of Mycobacterium tuberculosis infection. Biomed. Pharmacother. 91, 546–558. doi: 10.1016/j.biopha.2017.04.105 PubMed Abstract | CrossRef Full Text | Google Scholar Bajorath, J. (2012). Modeling of activity landscapes for drug discovery. Expert. Opin. Drug Discov. 7, 463–473. doi:

References AlMatar, M., AlMandeal, H., Var, I., Kayar, B., and Köksal, F. (2017). New drugs for the treatment of Mycobacterium tuberculosis infection. Biomed. Pharmacother. 91, 546–558. doi: 10.1016/j.biopha.2017.04.105 PubMed Abstract | CrossRef Full Text | Google Scholar Bajorath, J., ... [remainder of bibliography]

for the differences in the structural properties. In the classical QSAR studies, biological responses have been correlated with atomic, group, or molecular properties such as lipophilicity, polarizability, electronic, and steric properties (Hansch analysis) or with certain structural features (Free–Wilson analysis). However, in these techniques, one cannot ignore their limited utility for designing diverse functional new molecules due to the lack of consideration of the three-dimensional (3D) structures of the molecules. As a consequence, 3D-QSAR has emerged as a natural extension to the classical Hansch and Free–Wilson approaches that exploits the 3D properties of the ligands to predict their biological response by employing robust chemometric tools. The 3D-QSAR is a broad term encompassing all those QSAR methods that correlate macroscopic target properties with computed atom-based descriptors derived from the spatial representation of the molecular structures. These approaches have served as a valuable predictive tool in the design of pharmaceuticals and agrochemicals [1–3]. The prime goal of any 3D-QSAR method is to establish the relationship between biological activity and spatial properties of chemicals like steric, electrostatic, and lipophilic ones. The 3D-QSAR methodology is computationally more exhaustive and complex than 2D-QSAR approaches. Normally, it consists of several steps to acquire numerical descriptors from the compound structures. It is interesting to point out that some methods, independent of the alignment strategy, have also been developed with the progress of 3D-QSAR approaches [4]. One has to understand that the QSAR model is not a substitute for the experimental assays, although experimental techniques are also not free of inaccuracies. However, QSAR researchers are trying to develop a model that is as close as possible to the real one, and for this purpose, the 3D-QSAR techniques have to rely on some basic assumptions, which are illustrated here: • Binding of a drug molecule or ligand with the receptor is considered directly related to the biological response. Effects on second messengers or other signaling effects between receptor binding and experimentally observed response are not normally considered. • Molecular properties (physical, chemical, and biological) are encoded with a set of numbers or descriptors. • It is believed in general that compounds with common structures have comparable properties, and thus they have similar binding modes and accordingly equivalent biological activities and vice versa. • Structural properties leading to a biological response are usually determined by nonbonding forces, mainly steric and electrostatic ones. • Another important assumption is that the biological response is shown by the ligand itself, not by its metabolite product. • The lowest-energy conformation of the ligand is its bioactive conformation, which exerts binding effects. • The geometry of the receptor binding site is considered rigid, though there are a few exceptions. • The loss of translational and rotational degrees of freedom (entropy) upon binding is believed to follow a similar pattern for all these compounds. • The protein binding site is assumed to be the same for all of the studied ligands. • The major factors that contribute to the overall free energy of binding, like desolvation energy, temperature, diffusion, transport, pH, salt concentration, and plasma protein binding, are difficult to identify and thus are generally ignored. The 3D-QSAR methods can be classified based on a variety of criteria, as given in Table 8.1. Most commonly and successfully employed 3D-QSAR methods are discussed in the following sections of this chapter. Table 8.1 Categorization of 3D-QSAR techniques Basis of classification Type Examples of techniques Based on employed chemometric techniques Linear CoMFA, CoMSIA, AFMoC, GERM, CoMMA, SoMFA Nonlinear Compass Based on the alignment criterion Alignment-dependent CoMFA, CoMSIA, MSA, RSA, GERM, AFMoC, HIFA, VFA, MQSM Alignment-independent Compass, CoMMA, HQSAR, WHIM, GRIND, VolSurf, CoSA Based on the information employed to develop QSAR Ligand-based CoMFA, CoMSIA, MSA, RSA, Compass, GERM, CoMMA, SoMFA Receptor-based AFMoC, HIFA Comparative molecular field analysis (CoMFA) is a molecular field–based, alignment-dependent, ligand-based method developed by Cramer et al. [5], which helps in building the quantitative relationship of molecular structures and its response property. The method mostly focuses on ligand properties like steric and electrostatic ones, and the resulting favorable and unfavorable receptor–ligand interactions. As CoMFA is an alignment-dependent, descriptor-based method, all aligned ligands are placed in an energy grid, and by placing an appropriate probe at each lattice point, energy is calculated. The resultant energy calculated at each unit fraction corresponds to electrostatic (Coulombic) and steric (van der Waals) properties. These computed values serve as descriptors for model development. These descriptor values are then correlated with biological responses employing a robust linear regression method like partial least squares (PLS). The PLS results serve as an important signal to identify the favorable and unfavorable electrostatic and steric potential and also correlate it with biological responses. The formalism of the CoMFA methodology is described next: a. Structures of all molecules are drawn using any structure-drawing software. b. The bioactive conformation of each molecule is generated and energy minimization is carried out. c. All the molecules are superimposed or aligned using either manual or automated methods employed in the working software, in a manner defined by the supposed mode of interaction with the receptor. d. Thereafter, the overlaid compounds are positioned in the center of a lattice grid with a spacing of 2 Å. e. In the 3D space, the steric and electrostatic fields are calculated around the molecules with different probe groups positioned at all intersections of the lattice. Computation of the steric field uses the Lennard-Jones equation as follows: (8.1) In Eq. (8.1), $\varepsilon$ is the depth of the potential well, $\sigma$ is the finite distance at which the interparticle potential is zero, r is the distance between the particles, and rm is the distance at which the potential reaches its minimum. At rm, the potential function has the value $-\varepsilon$. The distances are given as $rm = 2^{1/6}\sigma$. Again, computation of electrostatic field follows the Coulombic interaction equation as follows: (8.2) where q1 and q2 denote point charges, r is the distance between charges, and $\varepsilon$ is the dielectric constant of the medium. f. The interaction energy or field values forming a pool of the descriptor/variable matrix are correlated with the biological response data employing the PLS technique, which identifies and extracts the quantitative influence of specific features of molecules on their activity. g. The results may be expressed as correlation equations with the number of latent variable terms, each of which is a linear combination of original independent lattice descriptors. h. For visual interpretation, the PLS output is illustrated in the form of interactive graphics consisting of colored contour plots of coefficients of the corresponding field variables at each lattice intersection, and showing the imperative favorable and unfavorable regions in the 3D space, which are closely associated with the biological activity. The CoMFA formalism is schematically illustrated in Figure 8.1. There are diverse factors that can control the complete performance of the constructed CoMFA model. These are described in the next sections. Like any 2D-QSAR method, one has to use precise activity data in order to create a good 3D-QSAR model. The following conditions should be fulfilled for maintaining the accuracy and appropriateness of the biological response data [3,6]: Representing the initial molecular structure is an important issue in 3D-QSAR analysis. This can be done by both experimental and computational approaches. A huge number of experimentally determined crystal structures are accessible in databases like the Cambridge Structural Database [7] and the Protein Data Bank [8]. The obtainable crystal structures present the benefit that some conformational information about the flexible molecule is included. Computationally, the 3D structures can be generated by three methods: Once the starting 3D molecular structures are generated, their geometries are refined by minimizing their conformational energies using following structure optimization techniques, including: • Molecular mechanics: It does not explicitly consider the electronic motion, so they are fast, accurate, and can be employed for large molecules like enzymes. • Quantum mechanics or ab initio: It takes into account the 3D distribution of electrons around the nuclei, and thus it is extremely precise. • Semiempirical: Semiempirical quantum chemical methods attempt to address two restrictions—namely, slow speed and low accuracy of quantum mechanical (e.g., Hartree–Fock) calculations by omitting certain integrals based on experimental data, such as ionization energies of atoms or dipole moments of molecules. Thus, semiempirical methods are very fast, applicable to large molecules, and may give precise results when applied to molecules that are similar to the molecules used for parameterization. Molecules to be used for semiempirical calculations may contain hundreds of atoms. Modern semiempirical models are based on the neglect of diatomic differential overlap (NDDO) methods like MNDO, AM1, PM3, and PDDG/PM3. The following conformational search methods can be implemented: • Systematic search (or grid search): It generates all probable conformations by systematically varying each of the torsion angles of a molecule by some increment, keeping the bond lengths and bond angles fixed. • Monte Carlo: It simulates dynamic behavior of a compound and generates the conformations by making random changes in its structure, calculating and comparing its energy with that of the previous conformation, and accepting the result if it is unique. • Random search: It generates a set of conformations by repetitively and arbitrarily changing either the Cartesian (x, y, z) or the internal (bond lengths, bond angles, and torsion/dihedral angles) coordinates of a starting geometry of the molecule under consideration. • Molecular dynamics: It employs Newton's second law of motion (force=mass×acceleration) to simulate the time-dependent movements and conformational changes in a molecular system, and results in a so-called trajectory showing how the positions and velocities of atoms in the molecular system vary with time. • Simulated annealing: It theoretically heats up the molecular system under consideration to high temperatures to overcome huge energy barriers, and after equilibrating there for some time using molecular dynamics, cools down the system slowly and gradually to obtain low-energy conformations according to the Boltzmann distribution. • Distance geometry algorithm: It generates a random set of coordinates by selecting random distances within each pair of upper and lower bounds to form constraints in a distance matrix, which are employed to create energetically feasible conformations of a set of molecules. • Genetic and evolutionary algorithms: It is based on the concept of biological evolution and initially creates a population of promising solutions to the problem. The solutions with the best fitness scores undergo crossovers and mutations over a time, and proliferate their good distinctiveness down the generations resulting in better solutions in the form of new conformers. The bioactive conformation defines a particular conformation of the molecule in which it is bound to the receptor. The intrinsic forces between the atoms in the molecule, as well as extrinsic forces between the molecule and its surrounding environment, considerably influence the bioactive conformation of the molecule [6]. Bioactive conformations of the compounds can be attained both by experimental and theoretical techniques. Experimental methods for creating bioactive conformations comprise the techniques described. The precise 3D structure of the macromolecules can be obtained by this method. Drug–receptor complexes generated by X-ray crystallography logically offer the exact information, but this method has several disadvantages: The 3D structural data is obtained in the solution and is a method of selection when the molecule cannot be crystallized through experimental ways, as in the case of the membrane-bound receptors or receptors, which have not yet been isolated due to stability, resolution, or other issues. The imperative features of this method are: