

# Statistical and Machine Learning Approaches to Predicting Protein-Ligand Interactions

Lucy J. Colwell

*Department of Chemistry, Cambridge University, Cambridge, UK*

---

## Abstract

Data driven computational approaches to predicting protein-ligand binding are currently achieving unprecedented levels of accuracy on held-out test data sets. Up until now, however, this has not led to corresponding breakthroughs in our ability to design novel ligands for protein targets of interest. This review summarizes the current state of the art in this field, emphasizing the recent development of deep neural networks for predicting protein-ligand binding. We explain the major technical challenges that have caused difficulty with predicting novel ligands, including the problems of sampling noise and the challenge of using benchmark datasets that are sufficiently unbiased that they allow the model to extrapolate to new regimes.

---

## Introduction

Molecular recognition is a fundamental requirement of biological systems. The interactions between proteins and small molecules are central to biology, allowing cells to sense their surroundings and respond appropriately. Estimates  
5 place the number of small molecules that can be synthesized at  $\approx 10^{60}$ , yet just a small fraction of potential protein-ligand interactions have been explored. Finding novel interactions is of great importance to drug discovery and basic biology. Given the enormity of the search space, computational approaches can narrow down the possibilities. However, despite three decades of computational  
10 effort, biochemical experiments are still essential to determine the efficacy of

ligand binding to a protein target [1, 2]. The results of computational analysis have been decidedly mixed: it is challenging to use even experimentally well-characterized ligand protein interactions to computationally design novel interactions [1, 2], much less explore the vast space of possibilities.

15 There are three increasingly demanding tasks in protein-ligand binding prediction: *virtual screening* predicts whether a ligand binds to a given target; *affinity prediction* predicts the binding affinity; and *pose prediction* identifies the molecular interactions causing binding to occur. In this review we focus on the first; the others have been reviewed elsewhere [3, 4]. Approaches to virtual  
20 screening can be categorised as physical or statistical. The idea of using first principles physical models to describe protein-ligand interactions is attractive, however timescale and computational resource constraints mean that simplified descriptions of features such as protein flexibility, and solvent are necessary. Even the most sophisticated docking algorithms cannot accurately reproduce  
25 large numbers of known interactions, much less predict new ones. Scoring functions can be empirical [5, 6, 7, 8, 9] or knowledge-based [10, 11, 12, 13, 14]. Significant expertise is required to encode physico-chemical interactions through the use of hand-tuned features and parameters, and can be highly specific to the system that they are designed for [15].

30 Recently, the use of high throughput methods to screen large libraries of proteins and small molecules and quantify their interactions has made it possible to correlate activity with representations of proteins and small molecules, to infer predictive models. Techniques from machine learning and artificial intelligence have been introduced, allowing both the parameters and the model to  
35 be learned from the wealth of experimental data available in databases such as ChEMBL [16, 17, 18, 19]. Increasingly publications are demonstrating that data driven approaches have the potential to make significant contributions to these problems [20\*\*, 21\*, 22\*, 23\*] [24, 25, 26, 27].

## Machine learning - potential and limitations

40 The aim of any machine learning or statistical approach is to identify patterns among training examples that can be used to make predictions outside the training set. An algorithm achieves this by mapping the training set representation to a space in which active and inactive ligands segregate - this mapping can be guided by physical models [13][23\*] but is more often learned directly from the  
45 data without addition of extensive physico-chemical knowledge [21\*, 22\*][26]. Once this mapping has been learned, it is hoped that the location of new examples in this space - clustered with training set actives or inactives - will accurately predict their activity. To compare the performance of different algorithms, a test set of ligands with known activity is held-out during the training  
50 process. An algorithm that can make accurate predictions for this unseen data is presumed to extrapolate well.

In particular, approaches that use deep neural networks (DNNs) have been shown to make predictions on held-out test sets with eye-opening levels of accuracy, exceeding 0.95 AUC (Area under the Receiver Operator Characteristic)  
55 on benchmark datasets [28][22\*]. However, the extent to which such results extrapolate is not yet clear - are they overfit to the training data [29, 1][30\*\*]? A number of studies posit that the test/train protocol is not as exacting as it might appear due to the non-uniformity of the distribution of ligands in chemical space. If training set actives are closer to test set actives than to training set  
60 inactives, by some metric that is not fully predictive of protein-ligand binding activity such as molecular weight, or number of ring systems, then the algorithm can appear to make accurate predictions for test set molecules without being able to extrapolate this predictive ability [31, 19, 32]. Machine learning performs best when abundant data drawn uniformly from the space of interest  
65 is available, but in this setting human chemists choose which molecules to work with, often based on clear similarities to known success stories [33, 34].

Definitively showing that these approaches generalise is perhaps the outstanding challenge facing this field today. In the search for novel pharmaceu-

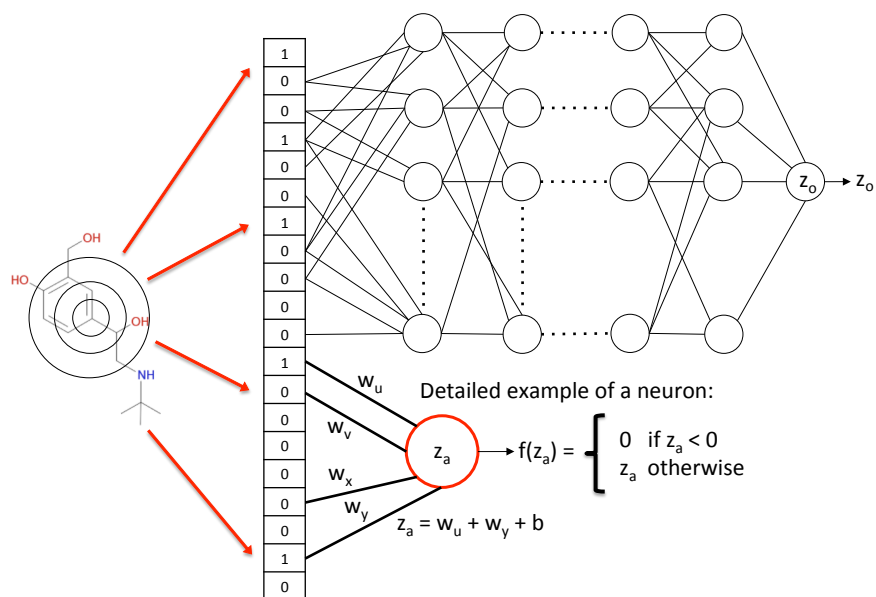


Figure 1: **Schematic of a deep neural network for predicting properties of small molecule ligands.** In this schematic, the input molecule is first encoded in a binary molecular fingerprint of fixed length. Convolutional architectures are a popular choice among practitioners and typically require that a fixed length representation be either used or learned from the molecular graph. The neural network simply consists of a very large number of simple functions (or neurons) of the form indicated in the schematic. Deep neural networks have large numbers of layers of neurons, resulting in a huge number of parameters that are learnt during training from the data, which is typically split into distinct training, cross-validation and test datasets.

70 ticals, the ability to predict the binding of ligands that are chemically distinct from those in the training data is highly valuable, but much more challenging for algorithms that are expert at identifying patterns among training set ligands. The goal of this article is to review statistical approaches to molecular recognition in the context of protein-ligand binding, focusing on recent results that exploit DNNs (see figure 1). Here we briefly outline the basic steps of machine  
 75 learning algorithm that predicts protein-ligand binding.

## Molecular Representation

There are almost as many choices for representation of the input data as there are for the machine learning algorithm employed [35, 36]. The simplest involve counting the numbers of different heavy atoms present in a ligand, together with  
80 other features such as hydrogen bond donors/acceptors, chiral centers and ring systems [19][30\*\*]. Some information about the chemical structure is retained by descriptors such as atom pairs or donor-acceptor pairs [37, 38] where each element has the form (atom type  $i$ ) - (distance in bonds) - (atom type  $j$ ). More information is encoded by chemical fingerprints, for example MACCS  
85 keys [39] and ECFP fingerprints [40]; fixed length binary descriptors which can be generated by the package RDKit [41]. Here, each non-hydrogen atom is used as a centre from which fragments are generated by extending radially from the centre along bonds to neighbouring atoms; the maximum radius considered  $N$  is encoded in the name as ECFP2 $N$ . A unique identifier is assigned to each  
90 fragment, and the set of identifiers for a molecule is mapped to a fixed length bit vector to yield the molecular fingerprint.

This abundance raises the question of which representation is most useful for different prediction tasks. Recently the suggestion has been made that it may be more effective to also learn the molecular representation itself, alongside  
95 the metric and corresponding embedding space used to distinguish active from inactive ligands [21\*, 22\*][27, 42]. However, counter-intuitively it has also been reported that the use of more complex molecular descriptors can result in little gain of predictive ability [43][30\*\*]

## Representation and Sampling Noise

100 One rationale for the finding that more complex representations can result in little improvement is noise due to finite sampling. The basic premise of any predictive algorithm is that similarities among known interaction partners can reveal the requirements of the binding site, and thus predict novel interactions. A straightforward approach is to compile the set of ligands known to bind to a

105 protein receptor of interest, and identify those features that show statistically  
significant enrichment among this set [44]. However, because there are only  
finitely many samples (i.e. known ligand binders), some features will be enriched  
purely by chance. This chance similarity increases with the number of variables,  
so representations that have more variables will lead to greater random similarity  
110 between features. For DNNs in particular, representations with thousands or  
even millions of features have recently been employed [24][21\*, 22\*]; although  
these algorithms have the ability to share information between targets it is  
still important that the level of chance similarity between small molecules is  
quantified and accounted for.

115 This phenomena has been carefully studied in the field of random matrix  
theory, which provides a null distribution that describes the similarity between  
samples (ligands) that can be expected by chance due to finite sampling as  
a function of the number of samples available, and the number of variables  
present in the ligand descriptor [45][46\*]. A simpler method for generating this  
120 null distribution simply involves computing the covariance matrices of multiple  
sets of  $n$  random ligands, where  $n$  is the sample size, using the same ligand  
descriptor for each set, and so obtaining the distribution of the largest entries  
that occur due to finite sampling noise. A similar approach can be taken for  
any measure of molecular similarity, defining a statistical null distribution with  
125 which to compare putative active molecules.

## Benchmark Datasets

The application of machine learning algorithms to problems such as protein-  
ligand binding is facilitated by the availability of large experimental datasets.  
For both algorithms that require 3D complexes and those that work with 2D  
130 ligand structures, a number of benchmark datasets have emerged that allow  
researchers to compare the performance of different algorithms. The PDB-  
bind benchmark [47] comprises 1300 diverse protein-ligand complexes with high  
quality structural and binding data. Another popular choice is the Directory

of Useful Decoys (DUD) [48] or Directory of Useful Decoys Enhanced (DUDE)  
135 benchmark [18], which have been carefully constructed to attempt to control  
for dataset biases that enable a machine learning algorithm to learn features  
that distinguish actives from inactives without providing information about  
the relevant molecular interaction. The construction of benchmark datasets  
that are unbiased, and so provide an accurate indication of how well an algo-  
140 rithm can extrapolate is an important research direction that requires atten-  
tion [49\*][50][30\*\*].

### Choice of Algorithm

Early linear regression methods were soon super-seeded by nonlinear models  
constructed using early versions of neural networks, among other techniques.  
145 Algorithms such as support vector machines [51, 52], Gaussian Processes [53, 54]  
and Random Forests (RF) [55, 15, 56, 43] have all received significant attention  
within the community. One advantage of these approaches is that many allow  
some level of interpretability, enabling the medicinal chemist to rationalise and  
evaluate the resulting model. This is more challenging for DNNs (see figure 1),  
150 which have recently been found to be highly successful on a variety of different  
tasks including computer vision [57] and speech recognition [58]. In the last few  
years there has been a burst of activity involving the construction of DNNs to  
predict protein-ligand binding [20\*\*, 21\*, 22\*] [24, 25, 26, 28, 59, 60, 61, 62].

A key breakthrough came in 2012 with the Merck sponsored Kaggle contest,  
155 in which the winning entry submitted by G. Dahl improved the mean squared  
Pearson correlation coefficient ( $R^2$ ) between predicted and observed binding  
activities over 15 proprietary Merck datasets from 0.42 (attained using RFs)  
to 0.49 [20\*\*]. Although a combination of DNNs [63\*], gradient boosting ma-  
chine [64] and Gaussian process regression [53] were used by the winning entry,  
160 Dahl suggests that the DNN was the main cause of the performance increase.  
A careful comparison of DNNs and RFs at protein-binding prediction over 15  
datasets found that the best  $R^2$  values over all parameter sets are obtained

using DNNs for five datasets, and using RFs for one, with mixed results for the remaining datasets. For refined parameter ranges, this increased to nine  
165 datasets for DNNs compared to one for RFs, leading the authors to conclude that DNNs have the potential to outperform RFs generally [20\*\*].

A crucial issue with DNNs is the set of algorithmic parameters affecting predictive performance that must be decided by the researcher, such as network size, choice of activation functions and use of dropout [65], which are distinct  
170 from parameters that are learned from the data. In general it is not computationally feasible to search the space of parameter values, in particular because different parameters are not independent of each other. Some analysis has been presented by varying one parameter at a time, and recording the  $R^2$  values of the resulting DNNs [20\*\*].

## 175 **Conclusion**

Advances in technology have led to enormous opportunity for using machine learning algorithms to learn models of protein ligand interactions. Machine learning algorithms are inherently flexible and with a sufficient amount of carefully curated data, the algorithms have the potential to find patterns  
180 that humans cannot find on their own. However, besides the issues already discussed with this review, other challenges must be overcome: the parameters learned from the data can lead to overfitting - the unintended specialisation of the trained model to the dataset used for training. For example, if the training and test datasets were identical, the algorithm would simply memorise all details of the training data, and achieve perfect performance on the test dataset.  
185 However, it would likely be unable to make accurate predictions for additional molecules that it had not previously seen before.

For chemical data, a key problem is to define what constitutes a distinct partitioning between test and training data. Simply requiring that the molecules  
190 are different from one another is not enough [31, 19, 32][30\*\*]. Machine learning algorithms are extremely effective at learning any feature that distinguishes

samples from different classes (such as binding and non-binding). This means that it is necessary to go further, and require that active molecules do not share any nuisance similarity or property that distinguishes them from inactive  
195 molecules, while not being the crucial molecular feature that enables binding. The situation is complicated by the fact that in most cases the crucial feature that enables binding is not known - indeed identifying these feature(s) is the ultimate goal of the exercise. However relatively simple devices such as synthetic datasets, in which binding is said to occur if particular logical combinations of  
200 features are present, or bias measures that quantify simple molecular similarities can be surprisingly effective.

Despite these challenges, this is an exciting time for building models of protein-ligand interactions. One could hope that a review written five years hence would be filled more with victories than warnings.

## 205 **References**

- [1] A. Peón, S. Naulaerts, P. J. Ballester, Predicting the reliability of drug-target interaction predictions with maximum coverage of target space, *Scientific Reports* 7 (2017).
- [2] P. C. Rathi, R. F. Ludlow, R. J. Hall, C. W. Murray, P. N. Mortenson, M. L. Verdonk, Predicting hot and warm spots for fragment binding, *Journal of Medicinal Chemistry* 60 (9) (2017) 4036–4046.  
210
- [3] R. Baron, J. A. McCammon, Molecular recognition and ligand association, *Annual review of physical chemistry* 64 (2013) 151–175.
- [4] J. D. Durrant, J. A. McCammon, Molecular dynamics simulations and drug discovery, *BMC biology* 9 (1) (2011) 71.  
215
- [5] M. D. Eldridge, C. W. Murray, T. R. Auton, G. V. Paolini, R. P. Mee, Empirical scoring functions: I. the development of a fast empirical scoring function to estimate the binding affinity of ligands in receptor complexes, *Journal of computer-aided molecular design* 11 (5) (1997) 425–445.

- 220 [6] H.-J. Böhm, The development of a simple empirical scoring function to estimate the binding constant for a protein-ligand complex of known three-dimensional structure, *Journal of computer-aided molecular design* 8 (3) (1994) 243–256.
- [7] R. Wang, L. Lai, S. Wang, Further development and validation of empirical  
225 scoring functions for structure-based binding affinity prediction, *Journal of computer-aided molecular design* 16 (1) (2002) 11–26.
- [8] R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E. H. Knoll, M. Shelley, J. K. Perry, et al., Glide: a new approach for rapid, accurate docking and scoring. 1. method  
230 and assessment of docking accuracy, *Journal of medicinal chemistry* 47 (7) (2004) 1739–1749.
- [9] O. Trott, A. J. Olson, Autodock vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multi-threading, *Journal of computational chemistry* 31 (2) (2010) 455–461.
- 235 [10] I. Muegge, Y. C. Martin, A general and fast scoring function for protein-ligand interactions: a simplified potential approach, *Journal of medicinal chemistry* 42 (5) (1999) 791–804.
- [11] H. Gohlke, M. Hendlich, G. Klebe, Knowledge-based scoring function to predict protein-ligand interactions, *Journal of molecular biology* 295 (2)  
240 (2000) 337–356.
- [12] H. Zhou, J. Skolnick, Goap: a generalized orientation-dependent, all-atom statistical potential for protein structure prediction, *Biophysical journal* 101 (8) (2011) 2043–2052.
- [13] M. L. Verdonk, R. F. Ludlow, I. Giangreco, P. C. Rathi, Protein–ligand  
245 informatics force field (pliff): Toward a fully knowledge driven force field for biomolecular interactions, *Journal of medicinal chemistry* 59 (14) (2016) 6891–6902.

- [14] G. Sliwoski, S. Kothiwale, J. Meiler, E. W. Lowe, Computational methods in drug discovery, *Pharmacological reviews* 66 (1) (2014) 334–395.
- 250 [15] P. J. Ballester, J. B. Mitchell, A machine learning approach to predicting protein–ligand binding affinity with applications to molecular docking, *Bioinformatics* 26 (9) (2010) 1169–1175.
- [16] A. Gaulton, A. Hersey, M. Nowotka, A. P. Bento, J. Chambers, D. Mendez, P. Mutowo, F. Atkinson, L. J. Bellis, E. Cibrián-Uhalte, et al., The chembl database in 2017, *Nucleic acids research* 45 (D1) (2016) D945–D954.
- 255 [17] M. K. Gilson, T. Liu, M. Baitaluk, G. Nicola, L. Hwang, J. Chong, Bindingdb in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology, *Nucleic acids research* 44 (D1) (2016) D1045–D1053.
- 260 [18] M. M. Mysinger, M. Carchia, J. J. Irwin, B. K. Shoichet, Directory of useful decoys, enhanced (dud-e): better ligands and decoys for better benchmarking, *Journal of medicinal chemistry* 55 (14) (2012) 6582–6594.
- [19] S. G. Rohrer, K. Baumann, Maximum unbiased validation (muv) data sets for virtual screening based on pubchem bioactivity data, *Journal of chemical information and modeling* 49 (2) (2009) 169–184.
- 265 [20] J. Ma, R. P. Sheridan, A. Liaw, G. E. Dahl, V. Svetnik, Deep neural nets as a method for quantitative structure–activity relationships, *Journal of chemical information and modeling* 55 (2) (2015) 263–274. \*\* Beyond the interesting results reported, that use deep neural nets (DNNs) to learn quantitative structure-activity relationships, this clearly written paper also uses careful analysis to shed light on the impact of DNN hyperparameters on predictive performance.
- 270 [21] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, R. P. Adams, Convolutional networks on graphs for

- 275 learning molecular fingerprints, in: *Advances in neural information processing systems*, 2015, pp. 2224–2232. \* Introduces an approach in which molecular descriptors are learned directly from the molecular graph. The authors report that these data-driven features are more interpretable, and display better performance on a number of prediction tasks.
- 280 [22] S. Kearnes, K. McCloskey, M. Berndl, V. Pande, P. Riley, Molecular graph convolutions: moving beyond fingerprints, *Journal of computer-aided molecular design* 30 (8) (2016) 595–608. \* The process of learning a molecular representation from the data is explained in this paper, which also explores different network architectures, and compares the resulting performance.
- 285 [23] A. P. Bartok, S. De, C. Poelking, N. Bernstein, J. Kermode, G. Csanyi, M. Ceriotti, Machine learning unifies the modelling of materials and molecules, arXiv preprint arXiv:1706.00179 (2017). \* The authors introduce a novel structure-based machine learning model that uses Gaussian process regression together with a local description of chemical environments to develop a unified framework that makes accurate predictions for problems that span from protein-ligand binding to the modelling of hard matter.
- 290 [24] T. Unterthiner, A. Mayr, G. Klambauer, M. Steijaert, J. K. Wegner, H. Ceulemans, S. Hochreiter, Deep learning as an opportunity in virtual screening, in: *Proceedings of the Deep Learning Workshop at NIPS*, 2014.
- 295 [25] I. Wallach, M. Dzamba, A. Heifets, Atomnet: a deep convolutional neural network for bioactivity prediction in structure-based drug discovery, arXiv preprint arXiv:1510.02855 (2015).
- 300 [26] G. B. Goh, C. Siegel, A. Vishnu, N. O. Hodas, N. Baker, Chemception: A deep neural network with minimal chemistry knowledge matches the performance of expert-developed qsar/qspr models, arXiv preprint arXiv:1706.06689 (2017).

- 305 [27] H. Altae-Tran, B. Ramsundar, A. S. Pappu, V. Pande, Low data drug discovery with one-shot learning, *ACS central science* 3 (4) (2017) 283–293.
- [28] B. Ramsundar, S. Kearnes, P. Riley, D. Webster, D. Konerding, V. Pande, Massively multitask networks for drug discovery, arXiv preprint arXiv:1502.02072 (2015).
- 310 [29] J. Gabel, J. Desaphy, D. Rognan, Beware of machine learning-based scoring functions on the danger of developing black boxes, *Journal of chemical information and modeling* 54 (10) (2014) 2807–2815.
- [30] I. Wallach, A. Heifets, Most ligand-based benchmarks measure overfitting rather than accuracy, arXiv preprint arXiv:1706.06619 (2017). \*\* Using  
315 a novel metric the authors demonstrate that many benchmark datasets are biased. This implies that machine learning algorithms trained on these datasets may not generalize beyond the domain of the training data as well as might be expected, given the reported performance on held out test data.
- 320 [31] M. L. Verdonk, V. Berdini, M. J. Hartshorn, W. T. Mooij, C. W. Murray, R. D. Taylor, P. Watson, Virtual screening using protein- ligand docking: avoiding artificial enrichment, *Journal of chemical information and computer sciences* 44 (3) (2004) 793–806.
- 325 [32] P. Ripphausen, A. M. Wassermann, J. Bajorath, Reprovis-db: a benchmark system for ligand-based virtual screening derived from reproducible prospective applications, *Journal of chemical information and modeling* 51 (10) (2011) 2467–2473.
- 330 [33] A. E. Cleves, A. N. Jain, Effects of inductive bias on computational evaluations of ligand-based modeling and on drug discovery, *Journal of computer-aided molecular design* 22 (3-4) (2008) 147–159.

- [34] A. N. Jain, A. E. Cleves, Does your model weigh the same as a duck?, *Journal of computer-aided molecular design* 26 (1) (2012) 57–67.
- [35] G. Maggiora, M. Vogt, D. Stumpfe, J. Bajorath, Molecular similarity in medicinal chemistry: miniperspective, *Journal of medicinal chemistry* 57 (8) (2013) 3186–3204.
- 335
- [36] A. Cereto-Massagué, M. J. Ojeda, C. Valls, M. Mulero, S. Garcia-Vallvé, G. Pujadas, Molecular fingerprint similarity search in virtual screening, *Methods* 71 (2015) 58–63.
- [37] R. E. Carhart, D. H. Smith, R. Venkataraghavan, Atom pairs as molecular features in structure-activity studies: definition and applications, *Journal of Chemical Information and Computer Sciences* 25 (2) (1985) 64–73.
- 340
- [38] S. K. Kearsley, S. Sallamack, E. M. Fluder, J. D. Andose, R. T. Mosley, R. P. Sheridan, Chemical similarity using physiochemical property descriptors, *Journal of Chemical Information and Computer Sciences* 36 (1) (1996) 118–127.
- 345
- [39] J. L. Durant, B. A. Leland, D. R. Henry, J. G. Nourse, Reoptimization of mdl keys for use in drug discovery, *Journal of chemical information and computer sciences* 42 (6) (2002) 1273–1280.
- [40] D. Rogers, M. Hahn, Extended-connectivity fingerprints, *Journal of chemical information and modeling* 50 (5) (2010) 742–754.
- 350
- [41] G. Landrum, Rdkit: Open-source cheminformatics, Online). <http://www.rdkit.org>. Accessed 3 (04) (2006) 2012.
- [42] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, G. E. Dahl, Neural message passing for quantum chemistry, *arXiv preprint arXiv:1704.01212* (2017).
- 355
- [43] P. J. Ballester, A. Schreyer, T. L. Blundell, Does a more precise chemical description of protein–ligand complexes lead to more accurate prediction

of binding affinity?, *Journal of chemical information and modeling* 54 (3) (2014) 944–955.

360 [44] R. Todeschini, V. Consonni, H. Xiang, J. Holliday, M. Buscema, P. Willett, Similarity coefficients for binary chemoinformatics data: overview and extended comparison using simulated and real data sets, *Journal of chemical information and modeling* 52 (11) (2012) 2884–2901.

[45] A. Edelman, Y. Wang, Random matrix theory and its innovative applica-  
365 tions, in: *Advances in Applied Mathematics, Modeling, and Computational Science*, Springer, 2013, pp. 91–116.

[46] A. A. Lee, M. P. Brenner, L. J. Colwell, Predicting protein–ligand affinity with a random matrix framework, *Proceedings of the National Academy of Sciences* 113 (48) (2016) 13564–13569. \* This paper applies random matrix  
370 theory to develop a novel approach that identifies the spurious enrichment of features among active or inactive ligands that can be caused by finite sampling effects.

[47] Z. Liu, Y. Li, L. Han, J. Li, J. Liu, Z. Zhao, W. Nie, Y. Liu, R. Wang, Pdb-  
wide collection of binding data: current status of the pddb database,  
375 *Bioinformatics* 31 (3) (2014) 405–412.

[48] N. Huang, B. K. Shoichet, J. J. Irwin, Benchmarking sets for molecular docking, *Journal of medicinal chemistry* 49 (23) (2006) 6789–6801.

[49] Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, V. Pande, Moleculenet: A benchmark for molecular machine  
380 learning, *arXiv preprint arXiv:1703.00564* (2017). \* Provides a comprehensive overview of data sets that can be used to train and test machine learning algorithms.

[50] N. Lagarde, J.-F. Zagury, M. Montes, Benchmarking data sets for the evaluation of virtual ligand screening methods: review and perspectives, *Journal of chemical information and modeling* 55 (7) (2015) 1297–1307.  
385

- [51] R. Burbidge, M. Trotter, B. Buxton, S. Holden, Drug design by machine learning: support vector machines for pharmaceutical data analysis, *Computers & chemistry* 26 (1) (2001) 5–14.
- [52] R. N. Jorissen, M. K. Gilson, Virtual screening of molecular databases using a support vector machine, *Journal of chemical information and modeling* 45 (3) (2005) 549–561.
- [53] F. R. Burden, Quantitative structure- activity relationship studies using gaussian processes, *Journal of chemical information and computer sciences* 41 (3) (2001) 830–835.
- [54] O. Obrezanova, G. Csányi, J. M. Gola, M. D. Segall, Gaussian processes: a method for automatic qsar modeling of adme properties, *Journal of chemical information and modeling* 47 (5) (2007) 1847–1857.
- [55] V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan, B. P. Feuston, Random forest: a classification and regression tool for compound classification and qsar modeling, *Journal of chemical information and computer sciences* 43 (6) (2003) 1947–1958.
- [56] D. Zilian, C. A. Sotriffer, Sfcscor rf: a random forest-based scoring function for improved affinity prediction of protein–ligand complexes, *Journal of chemical information and modeling* 53 (8) (2013) 1923–1933.
- [57] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [58] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al., Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups, *IEEE Signal Processing Magazine* 29 (6) (2012) 82–97.
- [59] S. Kearnes, B. Goldman, V. Pande, Modeling industrial admet data with multitask networks, *arXiv preprint arXiv:1606.08793* (2016).

- [60] A. Gonczarek, J. M. Tomczak, S. Zareba, J. Kaczmar, P. Dabrowski,  
415 M. J. Walczak, Learning deep architectures for interaction prediction in  
structure-based virtual screening, arXiv preprint arXiv:1610.07187 (2016).
- [61] M. Ragoza, J. Hochuli, E. Idrobo, J. Sunseri, D. R. Koes, Protein–ligand  
scoring with convolutional neural networks, *J. Chem. Inf. Model* 57 (4)  
(2017) 942–957.
- 420 [62] J. Gomes, B. Ramsundar, E. N. Feinberg, V. S. Pande, Atomic convolu-  
tional networks for predicting protein-ligand binding affinity, arXiv preprint  
arXiv:1703.10603 (2017).
- [63] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, MIT press, 2016.  
\* A comprehensive and clearly written introduction to machine learning  
425 algorithms and techniques.
- [64] V. Svetnik, T. Wang, C. Tong, A. Liaw, R. P. Sheridan, Q. Song, Boosting:  
An ensemble learning tool for compound classification and qsar modeling,  
*Journal of chemical information and modeling* 45 (3) (2005) 786–799.
- [65] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov,  
430 Dropout: a simple way to prevent neural networks from overfitting., *Journal*  
*of Machine Learning Research* 15 (1) (2014) 1929–1958.