
Dataset Bias in the Natural Sciences: A Case Study in Chemical Reaction Prediction and Synthesis Design

Ryan-Rhys Griffiths
Department of Physics
University of Cambridge
rrg27@cam.ac.uk

Philippe Schwaller
Department of Physics
University of Cambridge &
IBM Research Zurich
PHS@zurich.ibm.com

Alpha A. Lee
Department of Physics
University of Cambridge
aa144@cam.ac.uk

Abstract

Datasets in the Natural Sciences are often curated with the goal of aiding scientific understanding and hence may not always be in a form that facilitates the application of machine learning. In this paper, we identify three trends within the fields of chemical reaction prediction and synthesis design that require a change in direction. First, the manner in which reaction datasets are split into reactants and reagents encourages testing models in an unrealistically generous manner. Second, we highlight the prevalence of mislabelled data, and suggest that the focus should be on outlier removal rather than data fitting only. Lastly, we discuss the problem of reagent prediction, in addition to reactant prediction, in order to solve the full synthesis design problem, highlighting the mismatch between what machine learning solves and what a lab chemist would need. Our critiques are also relevant to the burgeoning field of using machine learning to accelerate progress in experimental Natural Sciences, where datasets are often split in a biased way, are highly noisy, and contextual variables that are not evident from the data strongly influence the outcome of experiments.

1 Introduction

Inventing new molecules through synthesis design is a central challenge for chemistry. Computers can process vast numbers of experimental reports and are able to accurately calculate the relative rates of competing reactions. It should therefore be the case that computers have the potential to produce more reliable synthetic routes than humans [1]. The synthesis design problem as well as the associated problems of reaction prediction and reaction planning are illustrated in Figure 1. The application of Machine learning to this problem has a history at NIPS [2] and has recently been demonstrated to be the state-of-the-art approach both in reaction prediction [3, 4, 5] and in synthesis design for reactants [6, 7, 8].

When married together, solutions to the reaction planning, reaction prediction and synthesis planning problems may be used to automatically propose routes to new molecules. Issues in dataset bias however are preventing the attainment of superhuman performance. In this paper, we identify three trends in the application of machine learning with respect to the design of synthetic pathways that may benefit from a change in direction. In section 2 we discuss reagent labelling in the reaction prediction problem and suggest an approach that may bring machine learning systems in line with industrial expectations. In section 3 we discuss the need for outlier detection in noisy reaction prediction datasets whilst in section 4 we highlight the fact that the prediction of reagents in addition to reactants in synthesis design is a key component in the design of a full synthetic route. We conclude by highlighting other areas in the Natural Sciences where dataset bias has been found to stymie progress, emphasizing the generality of the problem.

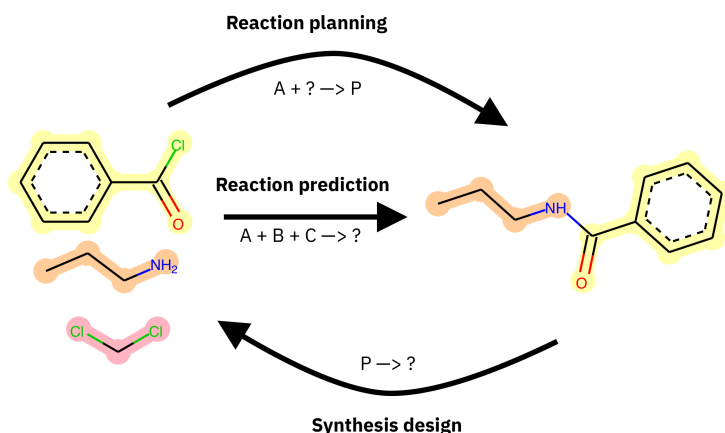


Figure 1: Designing Synthetic Pathways. The reaction planning problem involves finding a set of reagents to transform a given starting material into a given product. The reaction prediction problem concerns predicting the product given a set of reactants and the synthesis design problem involves working backwards from the product towards a set of reactants and reagents.

2 Reagent Labelling in Chemical Reaction Prediction

According to the IUPAC Compendium of Chemical Terminology a reagent is defined as “a test substance that is added to a system in order to bring about a reaction or to see whether a reaction occurs” [9]. In the chemical reaction prediction literature, the definition of a reagent is more precisely defined to be a compound that does not contribute atoms to the product.

For a synthetic organic chemist predicting the product of a new reaction, there is no way of knowing a priori which compounds in the reaction mixture will contribute atoms to the product because this is precisely the reaction prediction problem. Therefore, information about which chemicals are reactants and which chemicals are reagents implies prior knowledge about the product (see, e.g. [10], for algorithms that split reactants and reagents). As such, although the training set can be split into reactants and reagents, splitting the input of the test set into reactants and reagents makes the reaction prediction problem circular because the split can only be done with the answer known a priori. This dataset split is, however, routinely performed in current approaches [11, 5, 4, 12], where experiments are reported in which reagent labels have been provided to the model at test time. The distinction between providing and not providing reagent labels is illustrated in Figure 2.

[5] report results where the top-1 accuracy for the prediction model increases by 6% when reagents are labelled in the reaction prediction step. [11] explicitly label reagents using separate tokens. [4] input reagent information as a context vector to their model and [12] report results where improved performance is obtained with labelled reagents. Improved performance is not surprising in this case given that the space of possible products is narrowed through the exclusion of side reactions with the reagent. Since reagent labels are never available before a reaction is carried out, our recommendation would be for machine learning models to be benchmarked exclusively without reagent labelling. This has been done in [?] as well as the human comparison experiments of [3].

3 Noise in Chemical Reaction Data

Another difference between domains where machine learning has achieved significant progress, such as image recognition, and chemistry is the prevalence of label noise in chemical data. Image data benefits from negligible label noise, e.g. most images in MNIST and CIFAR-10 are correctly labelled (leading to the recent observation that algorithms which perfectly fit the training set can also achieve good generalisation error [13]). However, datasets in chemistry are often highly noisy. In chemical

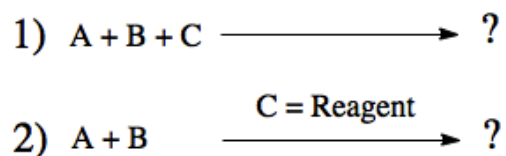


Figure 2: Reagent Labelling: 1) represents the realistic situation where no advance knowledge is available about which of the materials A, B or C acts as reactant/reagent. The question mark represents the product being predicted. 2) represents the output of many recent reaction prediction systems where C has been explicitly labelled as the reagent and is known to contribute no atoms to the product. In this instance, the system may be seen to be receiving some part of the ground truth label for the test-time prediction.

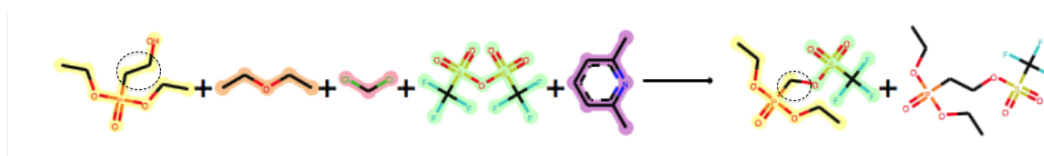


Figure 3: An example in the USPTO dataset of granted patent reactions where conservation of mass is violated. On comparing the circled region in the moiety on the extreme left-hand side of the figure with the moiety immediately to the right of the reaction arrow, we see that an extra carbon atom is present in the reactant that is not present in the product.

reaction prediction datasets such as the USPTO dataset of [14], there are numerous examples in which atoms are unduly present in the product, violating conservation of mass, as illustrated in Figure 3.

Although the community has, thus far, focused on improving model accuracy, the presence of a significant amount of noise in the dataset suggests that dataset cleaning and outlier removal is perhaps the key stumbling block to achieving superhuman performance. For example, “arrow pushing” – mapping out electron paths – is a sanity check tool in organic chemistry. The innovative preprocessing step of [4] models arrow pushing and prunes the data to consider only reactions that can be explained by linear electron topology. We argue that more effort should go into identifying and removing “impossible” chemical reactions from the dataset. We also speculate that a Bayesian approach which models aleatoric uncertainty – the inhomogeneous distribution of noise in the data – is an appropriate model to explore [15]. Specifically, the noise can come from (1) incorrect reporting of the structure in patents, (2) errors in transcribing the data into a digital format, as chemical structures extracted from patents are digitised using OCR technologies [16]. For both sources of noise, we argue that a reasonable conjecture is that the noise increases as a function of the chemical complexity of the reactants and reagents and is hence heteroscedastic.

4 Achieving the Automation of Synthesis Design

Although much progress has been made in the prediction of reactants given products [6, 7, 8], in order for the prediction to be actionable in the lab, one must predict both reactants and reagents given the products. This is because all chemical molecules can function as a reagent or a reactant depending on the chemical context and separating the two is an artificial construct as discussed above. The distinction between the problem solved by recent approaches and the full synthesis design problem is illustrated in Figure 4. More worryingly, the problem of reagent prediction that is seldom considered by the machine learning community is actually a challenging one – in fact, many Nobel prize-winning chemical innovation has centred on the discovery of new reagents [17, 18, 19]. As such, there is a gap between machine learning solutions and domain requirements.

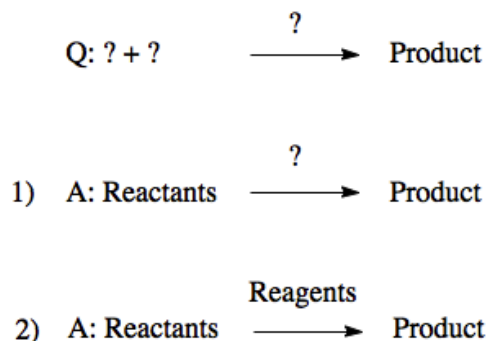


Figure 4: Towards Synthesis Design. The question is how to predict the reactants and reagents given the product. 1) Shows the problem currently being tackled, namely to only predict reactants given the product. The true goal however is 2) Predicting both reactants and reagents given the product in order to enable full automation of synthetic design.

5 Discussion

In this paper, we have focussed on dataset bias in chemical reaction prediction. Dataset bias however, is also prevalent in other areas; in ligand-based drug discovery, it has recently been shown that redundancy in the training and test sets can yield performance measures that do not accurately reflect industrial usage [20] whilst in the field of solubility prediction, literature values for diclofenac may differ by a factor of 100 [21] and so outlier detection should be an important preprocessing step.

In the field of molecule generation, the gap between machine learning system outputs and industrial requirements may benefit from the presence of more realistic benchmark objective functions. For example, the widely-used penalised logP objective [22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32] is likely to be too smooth a function of the molecular representation to resemble chemically interesting objective functions.

In accordance with the No Free Lunch Theorems [33, 34, 35], domain knowledge about the properties of industrially relevant objective functions may yield better heuristics for algorithm design. Objectives such as IC50 or binding affinity [36, 37, 38] are likely to be more interesting for drug discovery. To score highly for these properties, a certain level of specificity is required of the molecule and so one would expect the objective function to have sharp minima as opposed to broad minima as in the case of logP where similar molecules are expected to exhibit similar values of logP. Indeed positive work on correcting this trend of logP optimisation benchmarking has already been suggested by [39], who recommend a series of industrially relevant objectives.

6 Conclusion

In this paper we have highlighted three trends within the fields of chemical reaction prediction and synthesis design related to dataset presentation that need to be rethought: reagent labelling, outlier detection and reagent prediction. Our recommendations should hopefully redirect focus to problems that are of industrial relevance. We also discussed how problems associated with dataset presentation extend beyond the domain of reaction prediction. Perhaps outside the scope of Machine Learning but important to reaction prediction, we note that the chemistry literature is biased towards successful reactions and hence another dataset bias is the absence of data on reactions where the reactants do not react. Without non-reactions in the dataset, algorithms are biased towards predicting a chemical change in the reactants, whereas in reality it is not true that randomly mixing chemicals will always cause a chemical transformation. We suggest that a way to tackle this is through encouraging the adoption and sharing of Electronic Lab Notebooks in academic synthetic organic chemistry.

References

- [1] Jonathan M. Goodman. *Reaction Prediction and Synthesis Design*, chapter 4.2, pages 86–105. Wiley-Blackwell, 2018. ISBN 9783527806539. doi: 10.1002/9783527806539.ch4b. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/9783527806539.ch4b>.
- [2] Matthew A Kayala and Pierre F Baldi. A machine learning approach to predict chemical reactions. In *Advances in Neural Information Processing Systems*, pages 747–755, 2011.
- [3] Connor W. Coley, Wengong Jin, Luke Rogers, Timothy F. Jamison, Tommi S Jaakkola, William H. Green, Regina Barzilay, and Klavs F. Jensen. A Graph-Convolutional Neural Network Model for the Prediction of Chemical Reactivity. 10 2018. doi: 10.26434/chemrxiv.7163189.v1. URL https://chemrxiv.org/articles/A_Graph-Convolutional_Neural_Network_Model_for_the_Prediction_of_Chemical_Reactivity/7163189.
- [4] John Bradshaw, Matt J Kusner, Brooks Paige, Marwin HS Segler, and José Miguel Hernández-Lobato. Predicting electron paths. *arXiv preprint arXiv:1805.10970*, 2018.
- [5] Wengong Jin, Connor Coley, Regina Barzilay, and Tommi Jaakkola. Predicting organic reaction outcomes with weisfeiler-lehman network. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 2607–2616. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/6854-predicting-organic-reaction-outcomes-with-weisfeiler-lehman-network.pdf>.
- [6] Bowen Liu, Bharath Ramsundar, Prasad Kawthekar, Jade Shi, Joseph Gomes, Quang Luu Nguyen, Stephen Ho, Jack Sloane, Paul Wender, and Vijay Pande. Retrosynthetic reaction prediction using neural sequence-to-sequence models. *ACS central science*, 3(10):1103–1113, 2017.
- [7] Connor W Coley, Luke Rogers, William H Green, and Klavs F Jensen. Computer-assisted retrosynthesis based on molecular similarity. *ACS central science*, 3(12):1237–1245, 2017.
- [8] Marwin HS Segler, Mike Preuss, and Mark P Waller. Planning chemical syntheses with deep neural networks and symbolic ai. *Nature*, 555(7698):604, 2018.
- [9] A. D. McNaught and A. Wilkinson. *IUPAC. Compendium of Chemical Terminology, 2nd ed. (the "Gold Book")*. WileyBlackwell; 2nd Revised edition edition, August 1997. ISBN 978-0865426849.
- [10] Nadine Schneider, Nikolaus Stiefl, and Gregory A Landrum. What’s what: The (nearly) definitive guide to reaction role assignment. *Journal of chemical information and modeling*, 56(12):2336–2346, 2016.
- [11] Philippe Schwaller, Theophile Gaudin, David Lanyi, Costas Bekas, and Teodoro Laino. ?found in translation?: predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models. *Chemical science*, 9(28):6091–6098, 2018.
- [12] Anonymous. Graph transformation policy network for chemical reaction prediction. *Submitted to International Conference on Learning Representations*, 2019.
- [13] Yamini Bansal, Madhu Advani, David D Cox, and Andrew M Saxe. Minnorm training: an algorithm for training overcomplete deep neural networks. *arXiv preprint arXiv:1806.00730*, 2018.
- [14] Daniel Mark Lowe. *Extraction of chemical structures and reactions from the literature*. PhD thesis, University of Cambridge, 2012.

- [15] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017.
- [16] Nadine Schneider, Daniel M Lowe, Roger A Sayle, Michael A Tarselli, and Gregory A Landrum. Big data from pharmaceutical patents: a computational analysis of medicinal chemists’ bread and butter. *Journal of medicinal chemistry*, 59(9):4385–4402, 2016.
- [17] Georg Wittig and Werner Haag. Über triphenyl-phosphinmethylene als olefinbildende reagenzien (ii. mitteil. 1). *Chemische Berichte*, 88(11):1654–1666, 1955.
- [18] Norio Miyaura, Kinji Yamada, and Akira Suzuki. A new stereospecific cross-coupling by the palladium-catalyzed reaction of 1-alkenylboranes with 1-alkenyl or 1-alkynyl halides. *Tetrahedron Letters*, 20(36):3437–3440, 1979.
- [19] Tsutomu Katsuki and K Barry Sharpless. The first practical method for asymmetric epoxidation. *Journal of the American Chemical Society*, 102(18):5974–5976, 1980.
- [20] Izhar Wallach and Abraham Heifets. Most ligand-based classification benchmarks reward memorization rather than generalization. *Journal of chemical information and modeling*, 58(5):916–932, 2018.
- [21] Antonio Llinas, Jonathan C Burley, Karl J Box, Robert C Glen, and Jonathan M Goodman. Diclofenac solubility: independent determination of the intrinsic solubility of three crystal forms. *Journal of medicinal chemistry*, 50(5):979–983, 2007.
- [22] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- [23] Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. *arXiv preprint arXiv:1703.01925*, 2017.
- [24] Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained bayesian optimization for automatic chemical design. *arXiv preprint arXiv:1709.05501*, 2017.
- [25] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- [26] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for molecular graph generation. *arXiv preprint arXiv:1802.04364*, 2018.
- [27] Hiroshi Kajino. Molecular hypergraph grammar with its application to molecular optimization. *arXiv preprint arXiv:1809.02745*, 2018.
- [28] Zhenpeng Zhou, Steven Kearnes, Li Li, Richard N Zare, and Patrick Riley. Optimization of molecules via deep reinforcement learning. *arXiv preprint arXiv:1810.08678*, 2018.
- [29] Jiaxuan You, Bowen Liu, Rex Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. *arXiv preprint arXiv:1806.02473*, 2018.
- [30] Xiufeng Yang, Jinzhe Zhang, Kazuki Yoshizoe, Kei Terayama, and Koji Tsuda. Chemts: an efficient python library for de novo molecular generation. *Science and technology of advanced materials*, 18(1):972–976, 2017.
- [31] Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo drug design. *Science advances*, 4(7):eaap7885, 2018.
- [32] Daniil Polykovskiy, Alexander Zhebrak, Dmitry Vetrov, Yan Ivanenkov, Vladimir Aladinskiy, Polina Mamoshina, Marine Bozdaganyan, Alexander Aliper, Alex Zhavoronkov, and Artur Kadurin. Entangled conditional adversarial autoencoder for de novo drug discovery. *Molecular pharmaceutics*, 15(10):4398–4405, 2018.
- [33] Cullen Schaffer. A conservation law for generalization performance. In *Machine Learning Proceedings 1994*, pages 259–265. Elsevier, 1994.
- [34] David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- [35] David H Wolpert. The supervised learning no-free-lunch theorems. In *Soft computing and industry*, pages 25–42. Springer, 2002.

- [36] Shahar Harel and Kira Radinsky. Prototype-based compound discovery using deep generative models. *Molecular pharmaceutics*, 2018.
- [37] Tristan Aumentado-Armstrong. Latent molecular optimization for targeted therapeutic design. *arXiv preprint arXiv:1809.02032*, 2018.
- [38] Thomas Blaschke, Marcus Olivecrona, Ola Engkvist, Jürgen Bajorath, and Hongming Chen. Application of generative autoencoder in de novo molecular design. *Molecular informatics*, 37(1-2):1700123, 2018.
- [39] Daniel Neil, Marwin Segler, Laura Guasch, Mohamed Ahmed, Dean Plumbley, Matthew Sellwood, and Nathan Brown. Exploring deep recurrent models with reinforcement learning for molecule design. *Submitted to International Conference on Learning Representations*, 2018.