Article Review

"Best practices in machine learning for chemistry"

By Nongnuch Artrith, Keith T. Butler, François-Xavier Coudert, Seungwu Han, Olexandr Isayev, Anubhav Jain & Aron Walsh



Before we can dive into statistical learning, we must first create a solid foundation. The article "Best practices in Machine Learning in Chemistry" presents a collection of best practice standards and a matching checklist that can be used as a compass for navigating the challenging field of chemical machine learning. The researchers' goal was to maximize the validity and reproducibility of the conclusions and models derived from these best practices.

Some suggestions for using Machine Learning as a source of data:

The quality, quantity, and variety of accessible information have a significant impact on the efficacy of machine-learning models in chemistry.

Data can be gained by guided experiments or calculations, or they can be static and come from well-known chemical databases. The methodology used to generate the data or the setting in which a dataset was compiled are only two examples of the many elements that might lead to bias in data sources. Bias must be acknowledged and dealt with to produce clear and reliable models.

Additionally, databases in this discipline change over time, making version control systems and ongoing access to earlier dataset versions necessary to guarantee study reproducibility. It is crucial to list data sources, explain data selection methods, and give information about access times or version numbers.

Several suggestions for using Machine Learning as a model choice:

Machine learning provides a plethora of alternatives in the field of chemistry, ranging from conventional methods such as support vector machines to cutting-edge strategies like deep neural networks, particularly for graph-based chemical representations. Better outcomes aren't usually the product of more complicated models, though. For example, both a complicated neural network with over 13,000 parameters and a basic model with just two parameters worked. This emphasizes how crucial it is to select models carefully depending on the particular issue at hand.

When choosing a model, baseline comparisons—such as determining which class is more common or contrasting it with more basic models like the 1-nearest-neighbor approach—are essential, particularly when comparing the model to the most advanced models available.

It is also advised to provide a software implementation of the selected model to make training and testing with fresh data easier.

In summary, the field of Machine Learning in chemistry offers a range of options, but selecting the right model involves considering the complexity, conducting baseline comparisons, and ensuring accessibility through software implementations.

Some suggestions for using Machine Learning as a code and reproducibility:

The text analyzes the research reproducibility crisis, focusing on concerns such as selective reporting of good results, data dredging (p-hacking), and hypothesizing after results are known (HARKing). It emphasizes the significance of making research data and codes public to increase trust in science.

Explaining choices while training models are not adequate for reproducibility in Machine Learning. Hyperparameters and software versions are important considerations. As a result, the article suggests storing the whole code or workflow in a public repository for long-term replication and enhancement. A script or electronic notebook with all parameters should be provided as a bare minimum.

The most important message is the importance of transparency and accessibility in addressing the reproducibility challenge and bolstering scientific legitimacy.

In summary, the advancement of Machine Learning has brought about remarkable breakthroughs that have transformed our comprehension and management of chemical reactions. However, we must all adopt an open science approach to solve the issues of openness and reproducibility. This article's suggestions offer a road map for developing strong machine-learning models and making chemistry procedures easier to read.

Readers are recommended to click the linked <u>link</u> to read the complete article for a more thorough grasp of this topic.