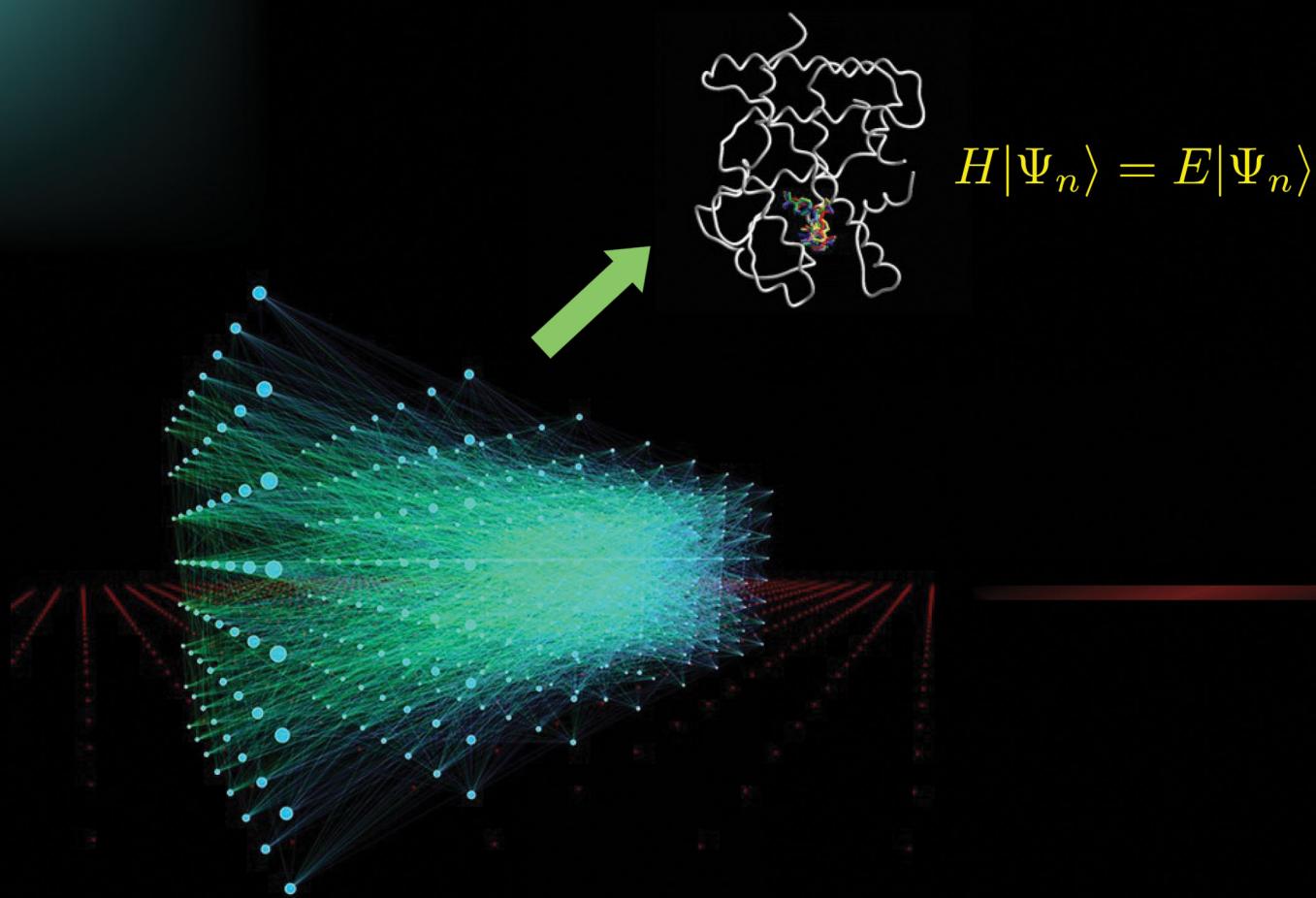


NJC

New Journal of Chemistry
rsc.li/njc

A journal for new directions in chemistry



ISSN 1144-0546

PERSPECTIVE

Richard Dybowski
Interpretable machine learning as a tool for scientific
discovery in chemistry



Cite this: *New J. Chem.*, 2020, **44**, 20914

Received 22nd May 2020,
Accepted 7th November 2020

DOI: 10.1039/d0nj02592e

rsc.li/njc

Interpretable machine learning as a tool for scientific discovery in chemistry

Richard Dybowski *

There has been an upsurge of interest in applying machine-learning (ML) techniques to chemistry, and a number of these applications have achieved impressive predictive accuracies; however, they have done so without providing any insight into what has been learnt from the training data. The interpretation of ML systems (*i.e.*, a statement of what an ML system has learnt from data) is still in its infancy, but interpretation can lead to scientific discovery, and examples of this are given in the areas of drug discovery and quantum chemistry. It is proposed that a research programme be designed that systematically compares the various model-agnostic and model-specific approaches to interpretable ML within a range of chemical scenarios.

1 AI and ML

Artificial intelligence (AI) is a large and complex subfield of computer science concerned with the development of algorithms that mimic (to some degree) human cognitive functions such as learning, image recognition and natural language processing.¹ It is believed that induction (generalising from finite examples) is an innate cognitive attribute,² and induction is the core interest of machine learning (ML),³ which has become a prominent part of AI.

Numerous techniques have been developed under the heading of ML including neural networks, support vector machines and random forests,⁴ but, in line with the concept of ‘statistical learning’,⁵ we will also regard all forms of statistical regression as being under the ML umbrella.

1.1 Deep learning

At the heart of neural-based computation is the concept of an artificial neural network.⁶ The term “deep neural networks” (DNNs) refers to neural networks that have several internal layers of neurons (Fig. 1), and their strength lies in their ability to make multiple non-linear transformations through these layers of neurons.⁷ In this process, increasingly complex and abstract features can be constructed by the addition of more layers and/or increasing the number of neurons per layer. Each layer can be thought of as performing an abstraction of the information held within the preceding layer, so that a sequence of layers provides a hierarchy of increasing abstraction. This can obviate the need for manual selection of input features.

St John's College, University of Cambridge, Cambridge CB2 1TP, UK.
E-mail: rd460@cam.ac.uk

1.2 Convolutional neural networks

Convolutional neural networks (CNNs) are a subclass of DNNs. CNNs⁸ are inspired by the Hubel–Wiesel model of the visual primary cortex⁹ in which complex images are built up from simple features. The standard architecture of a CNN consists of alternating convolutional and pooling layers. A convolutional layer preserves the relationship between values in a matrix by multiplying a submatrix of the layer with a matrix filter to produce a ‘feature map’. As a result, the network learns filters that activate when it detects some specific type of feature at some spatial position in the input.

A pooling layer abstracts the values of a feature map. This successive use of convolutional and pooling layers produces a hierarchy of abstracted features from an image that are invariant to translation, hence the successful use of CNNs for the recognition of images such as faces.

2 ML for chemistry

ML techniques, such as deep neural networks, have become an indispensable tool for a wide range of applications such as image

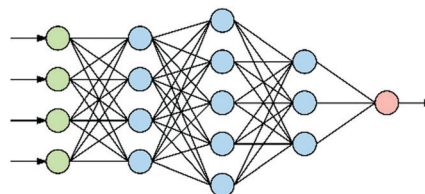


Fig. 1 An artificial neural network with three hidden layers (blue) between the input (green) and output (red) nodes: a deep neural network. Note that there can be more than 4 input nodes as well as more than one output node and many hidden layers.



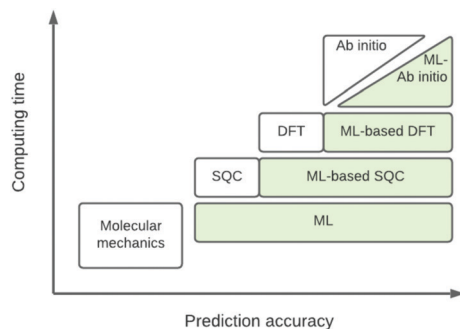


Fig. 3 Schematic representation of quantum chemical and ML approximations with respect to computational cost and accuracy, which generalises the literature.²⁴ DFT = density functional theory; SQC = semi-quantitative quantum chemistry.

3 Scientific discovery

The examples shown in the previous section focused on the use of ML for prediction, but what about scientific insight? Neural networks are so-called ‘black-box’ systems, meaning that the mapping of a vector of input values to a neural network’s output is too computationally complex for a human to comprehend; therefore, how can we, for example, discover what chemistry the AtomNet CNN has learnt from the vast amount of ligand–protein–interaction data used to train it for *in silico* drug screening?

The idea of using AI for scientific discovery is not new,²⁶ and there has recently been interest in using ML to provide scientific insights as well as making accurate predictions.²⁷

3.1 Interpretable ML

Currently, an important requirement of AI systems is not only the accuracy of the conclusions reached by the systems but also transparency as to how the conclusions were reached. Algorithms, particularly ML algorithms, are increasingly important to peoples’ lives, but they have caused a range of concerns revolving mainly around unfairness, discrimination and opacity. This has led to a ‘right to an explanation’ under the EU General Data Protection Regulation.

Interpretability is an ill-defined concept,²⁸ but it will suffice for us to use the definition that interpretable ML is the use of ML models for the extraction of relevant knowledge about domain relationships contained in data.²⁹ Consequently, interpretability refers to the extent to which a human expert can comprehend what an ML system has learnt from data; for example, ‘What is this ML system telling us?’ From an interpretation we have insight, and from insight we can hopefully make a scientific discovery.

There are several categories of interpretability in the context of ML. In the following list, $f(\mathbf{x})$ will be written as $f(\mathbf{x}; \hat{\Theta})$, where $\hat{\Theta}$ is the set of ML parameters estimated from training data.

(a) Observe the input–output behaviour of $f(\mathbf{x}; \hat{\Theta})$; for example, by observing how $f(\mathbf{x}; \hat{\Theta})$ varies as \mathbf{x} is varied.

(b) Inspect the values of parameters $\hat{\Theta}$ within the internal structure of $f(\mathbf{x}; \hat{\Theta})$. Here, $f(\mathbf{x}; \hat{\Theta})$ is either intrinsically interpretable or is interpretable by design. This allows mapping $\mathbf{x} \rightarrow f(\mathbf{x}; \hat{\Theta})$ to be understood by a series of steps going from input \mathbf{x} to output $f(\mathbf{x}; \hat{\Theta})$ that are comprehensible to a domain expert. This can be regarded as an ‘explanation’ of how $f(\mathbf{x}; \hat{\Theta})$ was derived from \mathbf{x} .

(c) Determine the prototypical value \mathbf{x} of for a given specific value f^* of $f(\mathbf{x}; \hat{\Theta})$. Conceptually, this can be regarded as the \mathbf{x} that maximises conditional probability $p(\mathbf{x}|f^*)$. \mathbf{x} need not have been previously encountered in a training set. An example of this approach is activation maximization.³⁰

A simple example of an intrinsically interpretable ML system is a linear regression model

$$\hat{\mathbb{E}}[y|x_1, \dots, x_n] = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

where β_0 , β_n are regression coefficients, and another is a decision tree induced from data.³¹

The potential of using interpretable ML for chemistry is starting to grow. For example, Bayesian neural networks have been optimised to predict the dissociation time of the unmethylated and tetramethylated 1,2-dioxetane molecules from only the initial nuclear geometries and velocities.³² Conceptual information was extracted from the large amount of data produced by simulations.

We now look at two other examples of interpretable ML: one from drug discovery; the other from quantum chemistry.

3.2 Drug discovery and interpretability

There are generally two approaches to providing interpretable ML: the model-agnostic and model-specific approaches (Fig. 4). Model-agnostic methods are, in principle, applicable to any black-box ML system $f(\mathbf{x})$, whereas the model-specific approach uses a domain-specific ML system, the structure of which is (at least partly) meaningful within the domain of interest.

There are two types of model-agnostic techniques. One method is the association-based technique in which associations are determined between inputs to system $f(\mathbf{x})$ and outputs from the system. One example of this are partial dependency plots,⁵ which create sets of ordered pairs $\{(\mathbf{x}_s^{(j)}, f(\mathbf{x}_s^{(j)}))\}$ where feature subset $\mathbf{x}_s \subset \mathbf{x}$. Another way to examine how $f(\mathbf{x})$ changes as x_j ($x_j \in \mathbf{x}$) changes is to use the ‘gradient input’ $x_i^* \partial f(\mathbf{x}) / \partial x_i$, where x_i^* is a particular value of x_i . However, an extension of this is to integrate the gradient along a path for x_i from observed value x_i^* to a baseline value x_i' . This is called an ‘integrated gradient’:

$$(x_i^* - x_i') \int_{\alpha=0}^1 \frac{\partial f(\tilde{\mathbf{x}})}{\partial \tilde{x}_i} \Big|_{\tilde{x}_i = x_i' + \alpha(x_i^* - x_i')} d\alpha.$$

Integrated gradients³³ determined the chemical substructures (toxicophores) that are important for differentiating toxic and non-toxic compounds. The relevant substructures identified in



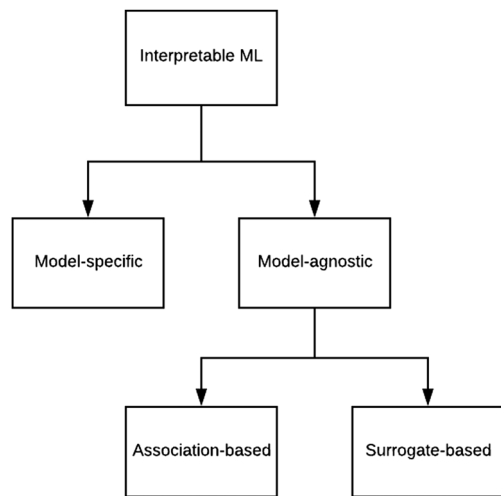


Fig. 4 Types of interpretable ML.

12 compounds randomly sampled from the Tox21 Challenge data set are shown in Fig. 5. The DNN consisted of four hidden layers, each with 2048 nodes. The molecular structures were encoded using ECFPs, the training and test sets had 12 060 and 647 examples, respectively, and the resulting AUC was 0.78.

An alternative to the detection of features relevant to a classification performed by a DNN is to start at an output node and work back to the input nodes. This is done with Layer-Wise Relevance Propagation (LRP),³⁴ which uses the network weights and the neural activations of a DNN to propagate the output (at layer M) back through the network up until the input layer (layer 1). The backward pass is a conservative relevance

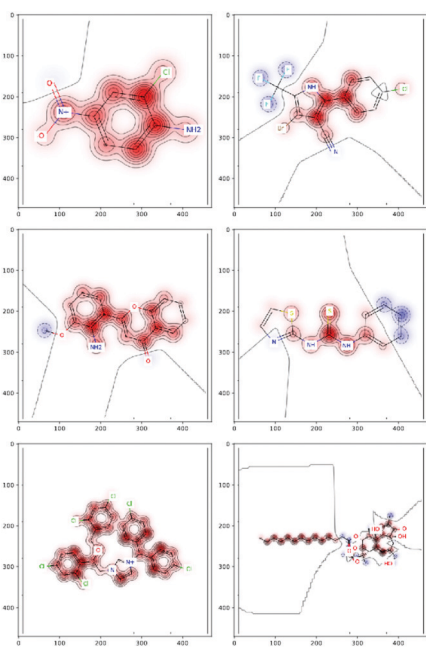


Fig. 5 Six randomly drawn Tox21 samples. Dark red indicates that these atoms are responsible for a positive classification, whereas dark blue atoms attribute to a negative classification.³³

redistribution procedure where those neurons in layer l ($1 \leq l < M$) that contribute the most to layer $l + 1$ receive the most 'relevance' from it. LRP has, so far, only been used to detect relevant features in pixel-based images and has not been used for the interpretation of chemistry-oriented ML systems. For example, rather than use fingerprints, such as ECPF, for molecular structure input, 2D molecular drawings have been used as inputs to a CNN (and achieved a predictive accuracy of AUC 0.766)³⁵ but LRP was not used for interpretation.

Graph convolutional neural networks (GCNNs) are a variant of CNNs that enable 3D graphs to be used as inputs. Consequently, if the nuclei and bonds of a compound are regarded as the vertices and edges of a 3D graph then 3D molecular structures can be considered as inputs to GCNNs. One approach is to initially slide convolutional filters over atom pairs to obtain atom-pair representations;³³ pooling is then used to produce simple substructure representations. These representations were then fed into the next convolutional layer. The predictive accuracy of the resulting GCNN was an AUC of 0.714. Interpretation was done by omitting the pooling steps and feeding the substructures directly into the fully connected network.

The other type of model-agnostic approach is the use of surrogates (Fig. 6); namely, using a function $\tilde{f}(x)$ that is an approximation of black box $f(x)$ but which is intrinsically interpretable. Examples of parsimonious intrinsically interpretable models include linear regression, logistic regression and decision trees. Such models can be either global or local. Given a set of vectors $\{x^{(1)}, \dots, x^{(n)}\}$ and that we wish to apply to $f(x)$, the global approach is to apply each of these vectors to the same surrogate model $\tilde{f}(x)$. In contrast, the local approach uses a different surrogate model $\tilde{f}_i(x^{(i)})$ for vector $x^{(i)}$. An example is the LIME technique,³⁶ which trains $\tilde{f}_i(x^{(i)})$ on data in the 'neighbourhood' of $x^{(i)}$, thereby providing interpretability specifically for the input-output pair $(x^{(i)}, \tilde{f}_i(x^{(i)}))$.

3.3 Quantum chemistry and interpretability

The above perspectives to interpreting a neural network focus on discovering associations between values at the input and output nodes: values at the hidden nodes are ignored. In contrast, another way of attempting to produce interpretable neural networks is to use internal nodes that are meaningful with respect to a domain of interest. This idea is not new and is at the heart of neuro-fuzzy systems³⁷ in which the interpretability of a neural network is provided by chains of inference *via* fuzzy logic.³⁸

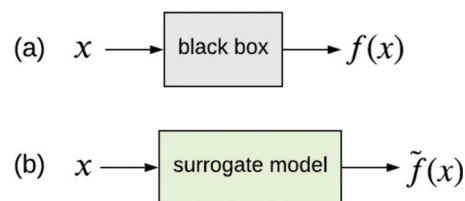


Fig. 6 (a) Black box trained from data $\{x, y\}$. (b) Surrogate model of the black box trained from the same data. $\tilde{f}(x)$ approximates $f(x)$.



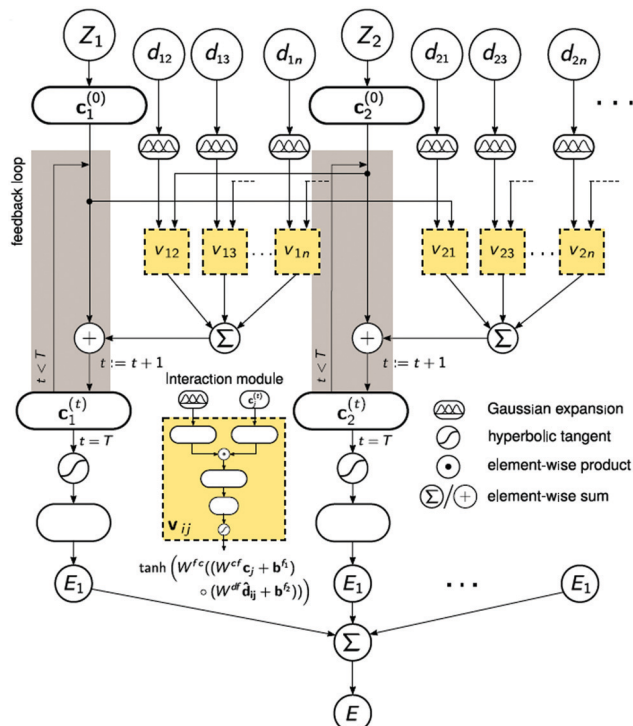


Fig. 7 The architecture of SchNet.³⁹ The iteration loop implements eqn (1), and the interaction module (a neural network) implements eqn (2).

Rather than resorting to fuzzy logic, neural networks such as SchNet (described below) are constructed by combining science-based subsystems in a plausible manner. This is an example of the model-specific approach to interpretable ML (Fig. 4).

A strategy for molecular energy E prediction³⁹ is to represent each atom i by a vector \mathbf{c}_i in B -dimensional space, and a deep tensor neural network (DTNN) called SchNet, shown in Fig. 7, repeatedly refines \mathbf{c}_i by pair-wise interaction between atoms i and j from an initial vector $\mathbf{c}_i^{(0)}$ for atom i to final vector $\mathbf{c}_i^{(T)}$:

$$\mathbf{c}_i^{(t+1)} = \mathbf{c}_i^{(t)} + \sum_{j \neq i} \mathbf{v}_{ij}, \quad (1)$$

Interaction term \mathbf{v}_{ij} reflects the influence of atom j at a distance \mathbf{d}_{ij} on atom i (the amount of overlap), and each refinement step aims to reduce these overlaps. For the interactions, the distances between atoms are expanded in a Gaussian basis.

Term \mathbf{v}_{ij} is obtained from atom vector \mathbf{c}_j and distance \mathbf{d}_{ij} using a feedforward neural network with a tanh activation function:

$$\mathbf{v}_{ij} = \tanh[W^{fc}(W^{cf}\mathbf{c}_j + \mathbf{b}^{f1}) \circ (W^{df}\mathbf{d}_{ij} + \mathbf{b}^{f2})] \quad (2)$$

where W^{cf} , \mathbf{b}^{f1} , W^{df} , \mathbf{b}^{f2} , and W^{fc} are the weight matrices and corresponding biases of atom representations, distances and resulting factors, respectively.

After T iterations, an energy contribution E_i for atom i is predicted for the final vector $\mathbf{c}_i^{(T)}$, and the total energy E is the sum of the predicted contributions E_i .

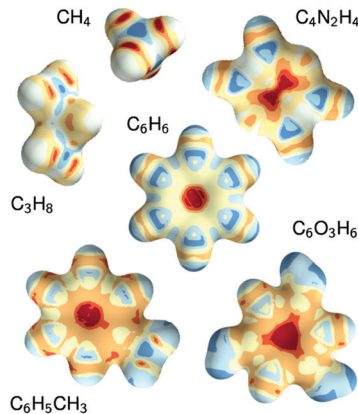


Fig. 8 Chemical potentials for methane, propane, pyrazine, benzene, toluene, and phloroglucinol determined from SchNet.³⁹

The DTNN, trained using stochastic gradient descent, achieved a mean absolute error of $1.0 \text{ kcal mol}^{-1}$ on the GDB datasets.

The constructive nature of the DTNN allows interpretation of how E is obtained, and the estimation of E allows energy isosurfaces to be constructed (Fig. 8).

Returning to the Schrödinger equation (in Dirac notation),

$$H|\Psi_n\rangle = E|\Psi_n\rangle$$

The Hartree–Fock method for molecular orbitals is to approximate a wave function $|\Psi\rangle$ for a molecular orbital as a linear combination of atomic orbitals:

$$|\Psi_n\rangle = \sum_{i=1}^N k_{n,i}|\phi_i\rangle$$

where $\{|\phi_i\rangle\}$ is a set of N basis functions, and $\{k_i\}$ are the associated coefficients. Function $|\phi_i\rangle$ can be an atom-centred Gaussian function. As a consequence, the electronic Schrödinger can be written in the matrix form

$$\mathbf{H}\mathbf{k}_m = \varepsilon_m\mathbf{S}\mathbf{k}_m$$

where the Hamiltonian matrix has elements

$$H_{i,j} = \langle\phi_i|H|\phi_j\rangle$$

and the overlap matrix has elements

$$S_{i,j} = \langle\phi_i|\phi_j\rangle$$

where $S_{i,j}$ measures the extent to which two basis functions overlap.

SchNOrb⁴⁰ was developed to predict \mathbf{H} and \mathbf{S} using ML. The first part of the structure of SchNOrb is identical to SchNet in that it starts from initial representations of atom types and positions, continues with the construction of representations of chemical environments of atoms and atom pairs (again identical to the method used in SchNet) but then uses these to predict energy E and Hamiltonian matrix \mathbf{H} , respectively.

The SchNet and SchNOrb systems illustrate how DNNs can be customized to specific scientific applications so that the DNN architecture promotes properties that are desirable in the



