# Machine Learning Methods in Geoscience

Gerard T. Schuster with labs by Yuqing Chen, Yongxiang Shi, Shihang Feng, Zhaolun Liu, Zongcai Feng, Yunsong Huang, Yan Yang and Tushar Gautam

March 14, 2022

# Contents

# Preface

This book presents the theory of machine learning (ML) algorithms and their applications to geoscience problems. More than half of the described algorithms fall under the class of neural network methods. Their description is at a level that can be understood by anyone with a modest background in linear algebra, calculus and probability. An elementary working knowledge of MATLAB is assumed and almost every chapter is accompanied by lab exercises to reinforce the ML principles. If you don't know MATLAB or Python, many of the labs use the CoLab, aka Colaboratory, notebook which is an easy-to-use product from Google that allows for the execution of Keras-based ML computer codes. It also provides for the free use of a cloud-based GPU processor[1] that can be used as soon as you login. There is no need for the troublesome installation of notebook software.

There are many definitions of *machine learning* (ML), but one of the best is attributed to Arthur Samuel[2]:

*Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.*

This definition doesn't mean we don't have to write the programs, it means that the ML algorithm finds the best model by learning from the data without relying on rules-based programming. For example, finding the best fit model $\mathbf{y} = f(\mathbf{x})$ by a typical least squares procedure typically assumes a linear mathematical model $f(\mathbf{x})$ that predicts the output $\mathbf{y}$ from the input $\mathbf{x}$. For example, $\mathbf{W}\mathbf{x} = \mathbf{y}$, where $\mathbf{W}$ is a matrix. A ML algorithm such as a neural network avoids this assumption by devising the best non-linear model learned from the data $\mathbf{y}$ for the given architecture.

**Neural Network Example.** How does the neural network learn the model[3] $f(\mathbf{x}) = \mathbf{y}$ from the data? Learning the model $f(\mathbf{x})$ is accomplished by using a large training set of data $(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ for $n \in \{1, 2, \ldots, N\}$, and adjusting the parameters of the model until $f(\mathbf{x}^{(n)})$ can accurately predict $\mathbf{y}^{(n)}$ for any training pair. Once the model is learned, then it can accurately predict the output $f(\mathbf{x}^{new}) = \mathbf{y}^{new}$ from a new input $\mathbf{x}^{new}$. This can be important

---

[1]Colab is a hosted Jupyter notebook service that requires no setup to use, while giving free access to computing resources including GPUs. See $https://colab.research.google.com/notebooks/intro.ipynb$

[2]This quote is often attributed to Arthur Samuel, but it cannot be found in Samuel (1959, 1967). A similar quote can be found at https://www.techemergence.com/what-is-machine-learning/. See http://infolab.stanford.edu/pub/voy/museum/samuel.html for a brief bio of this pioneering engineer.

[3]"Learning the model" is equivalent to computing the parameters of the non-linear function $f(\mathbf{x})$ that allow for an accurate prediction of $f(\mathbf{x}) = \mathbf{y}^{true}$ for all the training pairs $(\mathbf{x}, \mathbf{y}^{true})$.

for many applications, such as predicting the presence of underground contaminants $\mathbf{y}^{new}$ from processed geoelectrical and/or seismic data $\mathbf{x}^{new}$, accurately detecting the presence of breast cancer $\mathbf{y}^{new}$ recorded on a fusion of CAT scan and MRI images $\mathbf{x}^{new}$, or delineating the faults in $\mathbf{y}^{new}$ on the right-hand side of Figure 1 from the input seismic section $\mathbf{x}^{new}$ on the left-hand side.

The mathematical operations of a convolutional neural network model are represented by the block diagrams in Figure 1, where the input image $\mathbf{x}$ of a seismic section is decomposed into component images contained in the FM1 block (or shallowest layer). Here, a block represents a set of mathematical operations that produce a stack of 2D images known as feature maps (FMs). Each of the 32 feature maps in FM1 is computed by convolving a small filter ($3 \times 3$ filter in this example) with the input image, and then each pixel value in the FM is thresholded to be between 0 and 1 by a non-linear squashing function[4]. The $FM1$ feature maps are then subsampled by a factor of four (aka known as maxpooling), and the resulting images are filtered and thresholded by the FM2 $3 \times 3$ filters to get the 64 FMs for block 2. The subsampling acts as a low-pass spatial filter so that only coarser features of the seismic section are seen in, for example, FM3. In the end, the FMs represent the decomposition of a complicated object into high-wavenumber images on the left and low-wavenumber images on the far right. The FMs are weighted, upsampled, summed together and squashed to make a single decision, i.e. classification, at each pixel associated with the input image (Long et al., 2015). Does the pixel value represent a fault, then the answer is yes $y = 1$ for its class. If not, then the answer is no $y = 0$ for the output class. The displayed FMs are upsampled by a deconvnet procedure to transform them to the same size as the input image.

**Neural Networks and Deeper Connections.** Neural network models lack a rigorous theoretical foundation that explains their formidable number of commercial successes. This has motivated physicists and mathematicians to search for its mathematical and physical foundations. For example, they have recently discovered that successful CNN modeling can have many more tunable parameters than training data, i.e. equations of constraint, without suffering from overfitting. Arora et al. (2018) proves that a large number of FMs can be redundant and allows for simple compression of the model.

Another insight says that the CNN decomposition is similar to a dictionary-learning algorithm (Papayan et al., 2016, 2017a, 2017b; Sulam et al., 2018), where the different FMs can represent the multiscale components of a complicated object, such as its atoms or molecules. As an example, Figure 2 illustrates how two convolutional FMs, denoted as elementary "atoms" in the bottom row, are combined and filtered to create slightly more complex structures, i.e. molecules, at the second level. These "molecules" are then combined by a CNN model to create the global atom representing, in this case, a digit.

However, the basic operations of the neuron, i.e. atom, has much more complexity than previously thought, where the atom itself is composed of even finer structures. Work by Beniaguev et al. (2021) suggests that the cortical neuron mathematically performs the

---

[4]A squashing function $\sigma(z)$ *squashes* a wide range of input values $-\infty < z < \infty$ to be between a small range of output values such as $0 \leq \sigma(z) \leq 1$. The squashing operation sparsifies essential+extraneous information into its essential components (Papayan et al., 2016, 2017a, 2017b; Sulam et al., 2018; Chen et al., 2020b).

Figure 1: Block diagram of the convolutional neural network (CNN) architecture that decomposes the input picture of a seismic section into fine-scale feature maps (FM1 and FM2) and coarser-scale feature maps (FM3 and FM4). To provide a more interpretable image, the FMs are *deconvolved* for the display of interpretable images (Shi et al., 2018). The pixels in the far-right image are labeled $0 < y < 1$, where $y = 1$ indicates that the pixel is a fault and $0 \leq y \leq 1$ indicate the probability of being a fault. Typically, the shallow $3 \times 3$ filters that compute the 32 FM1 images from the input image detect the sharp edges in the input, and the filters associated with FM2 to FM3 detect larger scale features such as the general shape of the fault. These *successive decompositions derive the building blocks at successive levels from specific combinations of building blocks at the previous level* (Holland, 1995). Figure partially adapted from Dertat (2017) and Shi et al. (2018).

Figure 2: Elementary FMs in the bottom row are combined to form the molecules in the second row, which in turn are used to form the digit 6 in the top row. Figure from Sulam et al. (2018).

same functions as a deep neural network with 5-8 layers. This depth is associated with the neuron's dendritic branches that act as spatial pattern detectors.

Decomposing an input image into a sequence of simpler-to-more complicated parts is also used to describe the process of adaptive evolution illustrated in Figure 3. John Holland (1995) describes the fundamental set of procedures for adaptive evolution in nature:

*'If model making, broadly interpreted, encompasses most of scientific activity, then the search for building blocks becomes the technique for advancing that activity. At a fundamental level, we have the quarks of Gell-Mann. Quarks can be combined to yield nucleons, the building blocks at the next level. The process can be iterated, deriving the building blocks at successive levels from specific combinations of building blocks at the previous level. The result is the quark/ nucleon/ atom/ molecule/ organelle/cell/ underpins much of physical science.'.*

**Artificial Intelligence and CNN.** More generally, ML methods belong to the larger class of artificial intelligence[5] algorithms (see Figure 4), but the one that is attracting the most interest is that of deep learning with CNNs. Convolutional neural networks are now playing leading roles in the software for self-driving cars, language translation, speech recognition, computer vision, fusion and diagnosis of medical images, analysis of satellite and aerial

---

[5]https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai/

# Building Blocks and Recombination



A face can be described by stringing together the numbers that index its different component parts.

Figure 3: Neural networks is similar to John Holland's (1995) description of how nature evolves models: '...building blocks at successive levels from specific combinations of building blocks at the previous level.". The important features of the face are denoted by the *feature* rectangles on the left. Figure adapted from Holland (1995).

Figure 4: Artificial intelligence (AI) and the sub-classes of machine learning and deep learning, where AI involves a machine doing something that only a human would be able to do. Source of picture is Nvidia.

images for optimizing agricultural practices, and many other fields. Significantly, it is now being recognized as an important interpretation tool for analyzing lithology from seismic and well-log data, detecting faults in large 3D volumes of seismic data, and searching for patterns in earthquake data.

**Organization of the Book.** This book is divided into six main sections:

- Mathematical Background and Optimization Theory. Optimization theory and gradient descent methods for ML are reviewed, with an emphasis on the choice of step length methods. Readers with a familiarity of these algorithms can skip to the next sections.

- Supervised Learning Algorithms. Introduction to neural networks, fully-connected neural networks, applications to geoscience problems, and support vector machines.

- Convolutional Neural Networks. Basics of convolution and correlation, CNN, object identification, semantic segmentation, recurrent neural networks, Transformers and self-attention.

- Unsupervised Learning Algorithms: Autoencoders, convolutional sparse coding, principal component analysis, clustering methods, generative adversarial networks.

- Bayesian Analysis. Sampling, Bayes' theorem, and Gaussian mixture models.

- Seismic Inversion and ML. Physics-informed ML inversion, tomographic deconvolution, neural network LSM, wave equation inversion and neural networks, automatic differentiation.

Most chapters are accompanied by executable code (MATLAB and/or Keras with CoLab) that can be implemented by the reader with a modest level of programming knowledge. The understanding of the algorithms is deepened by parameter tuning and by examining the details of the code. Most of these codes have been written by my co-authors, without which this book would not have been written.

Case histories are used to demonstrate the practical use of ML in solving geoscience problems. These problems include fracture identification in photos, detection of rock cracks in drone photos, fossil and lithology detection in thin sections, prediction of permeability in rock samples, geochemical analysis, well-log analysis and identification of salt boundaries in seismic data, seismic arrival-time picking by a clustering method, least squares migration, velocity analysis, demultiple, noise reduction in seismic data and migration images, and some interpretation examples for seismic and radar data. The reader who diligently goes through the chapters and labs will have a thorough grounding in some of the fundamental ML methods used in geoscience.

**WWW Software Sites for Computational Labs.** The computational labs for the book *Machine Learning in Geoscience* are located at the following sites:

- *Machine Learning in Geoscience* Labs: http://csim.kaust.edu.sa/files/ErSE394/LAB1/

- *Machine Learning in Geoscience* Labs: http://utam.gg.utah.edu/books/ML/index.html

- *Machine Learning in Geoscience* Labs: http://seg.org/books/Schuster.ML/index.html

- *Machine Learning in Geoscience* Labs: repository.kaust.edu.sa/bitstream/handle/10754/674007/GeoscienceMachineLearning.html

- *Machine Learning in Geoscience* Labs: https://earth.utah.edu/books/Schuster.ML/index.html

The computational labs for the books *Seismic Interferometry* and *Seismic Inversion* are located at the following sites:

- *Seismic Interferometry* Labs: http://utam.gg.utah.edu/Inter.LAB1/

- *Seismic Inversion* Labs: repository.kaust.edu.sa/bitstream/handle/10754/674016/SeismicInversion.html

- *Seismic Inversion* Labs: https://earth.utah.edu/books/Schuster.SeismicInversion/index.html

# Acknowledgments