

Neural Networks in Imandra: Matrix Representation as a Verification Choice^{*}

Remi Desmartin¹, Grant Passmore², and Ekaterina Komendentskaya¹

¹ Heriot-Watt University, Edinburgh, UK

² Imandra Inc. Austin TX, USA lncs@springer.com

<http://www.springer.com/gp/computer-science/lncs>

³ ABC Institute, Rupert-Karls-University Heidelberg, Heidelberg, Germany
{abc,lncs}@uni-heidelberg.de

Abstract. The demand for formal verification tools for neural networks has increased as neural networks have been deployed in a growing number of safety-critical applications. Matrices are a data structure essential to formalising neural networks. Functional programming languages encourage diverse approaches to matrix definitions. This feature has already been successfully exploited in different applications. The question we ask is whether, and how, these ideas can be applied in neural network verification. A functional programming language Imandra combines the syntax of a functional programming language and the power of an automated theorem prover. Using these two key features of Imandra, we explore how different implementations of matrices can influence automation of neural network verification.

Keywords: Neural networks · Automated reasoning · Formal verification · Functional programming · Imandra.

1 Motivation

Neural network (NN) verification was pioneered by the SMT-solving [9, 10] and an abstract interpretation [1, 6, 18] communities. However, recently claims have been made that functional programming, too, can be valuable in this domain. There is a library [14] formalising small rational-valued neural networks in Coq. A more sizeable formalisation called MLCert [2] imports neural networks from Python, treats floating point numbers as bit vectors, and proves properties describing the generalisation bounds for the neural networks. An F^* formalisation [12] uses F^* reals and refinement types for proving robustness of networks trained in Python.

There are several options for defining neural networks in functional programming, ranging from defining neurons as record types [14] to treating them as functions with refinement types [12]. But we claim that two general considerations should be key to any NN formalisation choice of formalisation. Firstly, we

^{*} E.Komendentskaya acknowledges support of EPSRC grant EP/T026952/1.

must define neural networks as executable functions, because we want to take advantage of executing them in the functional language of choice. Secondly, a generic approach to layer definitions is needed, particularly when we implement complex neural network architectures, such as convolutional layers.

These two essential requirements dictate that neural networks are represented as matrices, and that a programmer makes choices about matrix formalisation. This extended abstract will explain these choices, and the consequences they imply, from the verification point of view. We use Imandra [15] to make these points, because Imandra is a functional programming language with tight integration of automated proving.

Imandra has been successful as a user-friendly and scalable tool in the Fin-Tech domain [16]. The secret of its success lies in combination of the best features of functional languages and interactive and automated theorem provers. Imandra’s logic is based on a pure, higher-order subset of OCaml, and functions written in Imandra are at the same time valid OCaml code that can be executed, or “*simulated*”. Imandra’s mode of interactive proof development is based on a typed, higher-order lifting of the *Boyer-Moore waterfall* [3] for automated induction, tightly integrated with novel techniques for SMT modulo recursive functions.

2 Matrices in Neural Network Formalisation

We will illustrate the functional approach to neural network formalisation and will introduce the syntax of the Imandra programming language [15] by means of an example. When we say we want to formalise neural networks as functions, essentially, we aim to be able to define a NN using just a line of code:

```
let cnn input =  
  layer_0 input >>= layer_1 >>= layer_2 >>= layer_3
```

where each `layer_i` is defined in a modular fashion.

To see that a functional approach to neural networks does not necessarily imply generic nature of the code, let us consider an example. A *perceptron*, also known as a *linear classifier*, classifies a given input vector $X = (x_1, \dots, x_m)$ into one of two classes c_1 or c_2 by computing a linear combination of the input vector with a vector of synaptic weights (w_0, w_1, \dots, w_m) , in which w_0 is often called an *intercept* or *bias*: $f(X) = \sum_{i=1}^m w_i x_i + w_0$. If the result is positive, it classifies the input as c_1 and if negative as c_2 . It effectively divides the input space along a hyperplane defined by $\sum_{i=1}^m w_i x_i + w_0 = 0$.

In most classification problems, classes are not linearly separated. To handle such problems, we can apply a non-linear function a called an *activation function* to the linear combination of weights and inputs. The resulting definition of a perceptron f is:

$$f(X) = a \left(\sum_{i=1}^m w_i x_i + w_0 \right) \quad (1)$$

Let us start with a naive prototype of perceptron in Imandra. The Iris data set is a “Hello World” example in data mining; it represents 3 kinds of Iris flowers using 4 selected features. In Imandra, inputs can be represented as a data type:

```
type iris_input = {
  sepal_len: real;
  sepal_width: real;
  petal_len: real;
  petal_width: real;}
```

And we define a perceptron as a function:

```
let layer_0 (w0, w1, w2, w3, w4) (x1, x2, x3, x4) =
  relu (w0 +. w1 *. x1 +. w2 *. x2 +. w3 *. x3 +. w4 *. x4)
```

where `*.` and `+.` are *times* and *plus* defined on reals. Note the use of the `relu` activation function, which returns 0 for all negative inputs and acts as the identity function otherwise.

Already in this simple example, one perceptron is not sufficient, as we must map its output to three classes. We use the usual machine learning literature trick and define a further layer of 3 neurons, each representing one class. Each of these neurons is itself a perceptron, with one incoming weight and one bias. This gives us:

```
let layer_1 (w1, b1, w2, b2, w3, b3) f1 =
  let o1 = w1 *. f1 +. b1 in
  let o2 = w2 *. f1 +. b2 in
  let o3 = w3 *. f1 +. b3 in
  (o1, o2, o3)

let process_iris_output (c0, c1, c2) =
  if (c0 >=. c1) && (c0 >=. c2) then "setosa"
  else if (c1 >=. c0) && (c1 >=. c2) then "versicolor"
  else "virginica"
```

The second function maps the output of the three neurons to the three specified classes. This post-processing stage often takes a form of an *argmax* or *softmax* function, which we omit.

And thus the resulting function that defines our neural network model is:

```
let model input = process_iris_input input
  |> layer_0 weights_0 |> layer_1 weights_1 |>
    process_iris_output
```

Although our naive formalisation has some features that we desired from the start, i.e. it defines a neural network as a composition of functions, it is too inflexible to work with arbitrary compositions of layers. In neural networks with

hundreds of weights in every layer this manual approach will quickly become infeasible (as well as error prone). So, let us generalise this attempt from the level of individual neurons to the level of matrix operations.

The composition of many perceptrons is often called a *multi-layer perceptron* (*MLP*). An MLP consists of an input vector (also called input layer in the literature), multiple hidden layers and an output layer, each layer made of perceptrons with weighted connections to the previous layers' outputs. The weight and biases of all the neurons in a layer can be represented by two matrices denoted by W and B . By adapting equation 1 to this matrix notation, a layer's output L can be defined as:

$$L(X) = a(X \cdot W + B) \quad (2)$$

where the operator \cdot denotes the dot product between X and each row of W , X is the layer's input and a is the activation function shared by all nodes in a layer. As the dot product multiplies pointwise all inputs by all weights, such layers are often called *fully-connected*.

By denoting a_k, W_k, B_k — the activation function, weights and biases of the k th layer respectively, an MLP F with L layers is traditionally defined as:

$$F(X) = a_L[B_L + W_L(a_{L-1}(B_{L-1} + W_{L-1}(\dots(a_1(B_1 + W_1 \cdot X)))))] \quad (3)$$

At this stage, we are firmly committed to using matrices and matrix operations. And we have two key choices:

1. to represent matrices as lists of lists (and take advantage of the inductive data type `List`),
2. define matrices as functions from pairs to matrix elements,
3. or take advantage of record types, and define matrices as records.

The first choice was taken in [8] (in the context of dependent types in Coq), in [12] (in the context of refinement types of F*) and in [7] (for sparse matrix encodings in Haskell). The difference between the first and second approaches was discussed in [20] (in Agda, but with no neural network application in mind). The third method was taken in [14] using Coq (though records there were used to encode individual neurons).

In the next three sections, we will systematise these three approaches using the same formalism and the same language, in order to understand the influence they make on neural network verification.

3 Matrices as Lists of Lists

We start with re-using Imandra's `List` library. Lists are defined as inductive data structures.

Remi, please give the full definition in Imandra's syntax?

We start with defining vectors as lists, and matrices as lists of vectors.

```

type 'a vector = 'a list
type 'a matrix = 'a vector list

```

It is possible to extend this formalisation by using dependent [8] or refinement [12] types to check the matrix size. But in Imandra this facility is not directly available, and we will need to use exceptions (monadic operations) to check the matrix sizes.

As there is no built-in type available for matrices equivalent to `List` for vectors, the `Matrix` module implements a number of functions for basic operations needed throughout the implementation. For instance, `map2` takes as inputs a function f and two matrices A and B of the same dimensions and outputs a new matrix C where each element $c_{i,j}$ is the result of $f(a_{i,j}, b_{i,j})$:

```

let rec map2 (f: 'a -> 'b -> 'c) (x: 'a matrix) (y: 'b matrix)
  = match x with
  | [] -> (match y with
    | [] -> Ok []
    | y::ys -> Error "map2: invalid list length.")
  | x::xs -> match y with
    | [] -> Error "map2: invalid list length."
    | y::ys -> let hd = map2 f x y in
      let tl = map2 f xs ys in
      lift2 cons hd tl

```

This implementation allows us to define other useful functions in a concise way. For instance, the dot-product of two matrices, or the L_0 distance between two matrices are defined as:

```

let dot_product (a:real matrix) (b:real matrix) =
  let c = map2 ( *. ) a b in
  map sum c

```

A fully connected layer is then defined as a function `fc` that takes as parameters an activation function, a 2-dimensional matrix of layer's weights and an input vector:

```

let activation f w i = (* activation func., weights, input *)
let linear_combination m1 m2 = if (length m1) <> (length m2)
  then Error "invalid dimensions"
  else map sum (Vec.map2 ( *. ) m1 m2) in
let i' = 1::i in (* prepend 1. for bias *)
let z = linear_combination w i' in
map f z

let rec fc f (weights:real matrix) (input:real vector) =
  match weights with
  | [] -> Ok []

```

```
| w::ws -> lift2 cons (activation f w input) (fc f ws input)
```

Listing 1.1: Fully connected layer implementation

Note that each row of the weights matrix represents the weights for one of the layer’s nodes. The bias for each node is the first value of the weights vector, and 1 is prepended to the input vector when computing the linear combination of weights and input to account for that.

It is now easy to see that our desired modular approach to composing layers works as stated. We may define the layers using the syntax: `let layer_i = fc a weights`, where `i` stands for 0,1,2,3, and `a` stands for any chosen activation function.

Although natural, this formalisation of layers and networks suffers from two problems. Firstly, it lacks the matrix dimension checks that were readily provided via refinement types in [12]. This is because Imandra is based on a computational fragment of HOL, and has no refinement or dependent types. To mitigate this, the library we present performs explicit dimension checking via a `result` monad, which clutters the code and adds additional computational checks. Secondly, the matrix definition via the list datatypes makes verification of neural networks very inefficient. This general effect has been already reported in [12], but it may be instructive to look into the problem from the Imandra perspective.

Robustness of neural networks [4] is best amenable to proofs by arithmetic manipulation. This explains the interest of the SMT-solving community in the topic, which started with using Z3 directly [9], and has resulted in highly efficient SMT solvers specialised on robustness proofs for neural networks [10,11]. Imandra’s waterfall method [15] defines a default flow for the proof search, which starts with unrolling inductive definitions, simplification and rewriting. As a result, proofs of neural network robustness or proofs as in the ACAS Xu challenge, which should not rely on the matrix size induction, stall in Imandra.

There is another mode of proofs available in Imandra: `blast`, a tactic for SAT-based symbolic execution modulo higher-order recursive functions. Blast is an internal custom SAT/SMT solver that can be called explicitly with the appropriate tactic. However, `blast` currently does not support real arithmetic. This requires us to *quantize* the neural networks we use (i.e. convert them to integer weights) and results in a *quantised NN library* [5]. However, even with quantisation and the use of Blast, Imandra fails to scale to the Acas Xu challenge, let alone neural networks used in computer vision.

However, there is a silver lining of this method of matrix formalisation. When we formulate verification properties that genuinely require induction, formalisation of matrices as lists does result in more natural, and easily automatable proofs. For example, De Maria et al. [14] formalise in Coq “*neuronal archetypes*” for biological neurons. Each archetype is a specialised kind of perceptron, in which additional functions are added to amplify or inhibit the perceptron’s outputs. It is out of scope of this paper to formalise the neuronal archetypes in Imandra, but we take methodological insight from [14]. In particular, [14] shows that there are natural higher-order properties that one may want to verify.

To make a direct comparison, modern neural network verifiers [10, 18] deal with verification tasks of the form “given a trained neural network f , and a property P_1 on its inputs, verify that a property P_2 holds for f ’s outputs”. However, the formalisation in [14] considers properties of the form “any neural network f that satisfies a property Q_1 , also satisfies a property Q_2 .” Unsurprisingly, the former kind of properties can be resolved by simplification and arithmetic, whereas the latter kind requires induction on the structure of f (as well as possibly nested induction on parameters of Q_1).

Another distinguishing consequence of this approach is that it is orthogonal to the community competition for scaling proofs to large networks: usually the property Q_1 does not restrict the size of neural networks, but rather points to their structural properties. Thus, implicitly we quantify over neural networks of any size.

To emulate a property *à la* de Maria et al., in [5] we defined a general network monotonicity property: *any fully connected network with positive weights is monotone, in the sense that, given increasing positive inputs, its outputs will also increase*. There has been some interest in monotone networks in the literature [17, 19]. Our experiments in [5] show that Imandra can prove such properties by induction on the networks’ structure almost automatically (with the help of a handful of auxiliary lemmas). And the proofs easily go through for both quantised and real-valued neural networks.

Remi, please check it is true? Anything else needs to be said here?

4 Matrices as Functions

5 Matrices as Records

6 Related Work and Conclusions

6.1 using Sized Lists or Vectors

Grant et al. ([7]) proposes 4 sparse list-based matrix implementations. They use an array-as-trees representation which allows to optimise for sparse arrays (subtrees where all the leaves are 0 are replaced by a 0-leaf).

Binary trees and lists of row-fragments: binary tree array of sparse Vectors defined as `[(Int, [Double])]`

A generalised envelope scheme: matrix is cut up in sections A quadtree scheme: Triangular matrix is split up in 2 triangular and a rectangular one. A standard quadtree structure is used for the rectangular matrix.

A Two-copy list of row-segments scheme: list of row-segments and list of column-segments in order to iterate over columns easily. Con: 2x more space is used; can be used to improve the 2 first previous methods (quadtree already bidimensional)

Pros of sparse list-based matrix representation: optimised for sparse matrices. Optimised for the specific operation considered in the paper (solving of linear systems of equations using a Cholesky scheme)

6.2 Refined Types

Coq/Mathcomp/SSReflex: the size of the matrix is defined as a refinement type and used to check that matrix have the appropriate size in multiplications Heras et al [8] discuss implementation; correctness is checked at compile-time.

This is also used in Starchild/Lazuli [13].

6.3 Arrays

Grant et al. [7] mentions using Haskell mutable arrays which are implemented using monadic operations. They stress that using a mutable array allows for modifying the array in place (thus saving memory), but it introduces "extra programming difficulties"; i.e. the use of monads makes the code less clear (as is the case in 6.4).

In our case, since IML is pure, there is no access to mutable arrays.

6.4 Monadic Operations to Check Size

Allows to check for valid matrix sizes when no refinement types are available. Intuitive first representation.

Drawback: introduces a lot of pattern matching, so a lot of split cases which increases the size of the program to check exponentially

Drawback: no optimisation for sparse matrices

6.5 Matrix as Functions

Imandra implementation; matrices are defined as total functions mapping indices to values. Theory of uninterpreted functions.

Woods [20] discusses the benefits of implementing matrices as functions in Agda.

6.6 Matrix as Maps

References

1. Ayers, E.W., Eiras, F., Hawasly, M., Whiteside, I.: PaRoT: A practical framework for robust deep neural network training. In: NASA Formal Methods - 12th International Symposium, NFM 2020, Moffett Field, CA, USA, May 11-15, 2020, Proceedings. LNCS, vol. 12229, pp. 63–84. Springer (2020)
2. Bagnall, A., Stewart, G.: Certifying true error: Machine learning in Coq with verified generalisation guarantees. AAAI (2019)
3. Boyer, R.S., Moore, J.S.: A Computational Logic. ACM Monograph Series. Academic Press, New York (1979)
4. Casadio, M., Komendantskaya, E., Daggett, M.L., Kokke, W., Katz, G., Amir, G., Refaeli, I.: Neural network robustness as a verification property: A principled case study. In: Computer Aided Verification (CAV 2022). Lecture Notes in Computer Science, Springer (2022)

5. Desmartin, R., Passmore, G., Komendantskaya, E., Daggett, M.L.: CNN library in Imandra. <https://github.com/aisec-private/ImandraNN> (2022)
6. Gehr, T., Mirman, M., Drachler-Cohen, D., Tsankov, P., Chaudhuri, S., Vechev, M.T.: AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation. In: S&P (2018)
7. Grant, P.W., Sharp, J.A., Webster, M.F., Zhang, X.: Sparse matrix representations in a functional language. *Journal of Functional Programming* **6**(1), 143–170 (Jan 1996). <https://doi.org/10.1017/S09567968000160X>, <https://www.cambridge.org/core/journals/journal-of-functional-programming/article/sparse-matrix-representations-in-a-functional-language/669431E9C12EDC16F02603D833FAC31B>, publisher: Cambridge University Press
8. Heras, J., Poza, M., Dénès, M., Rideau, L.: Incidence Simplicial Matrices Formalized in Coq/SSReflect. In: Davenport, J.H., Farmer, W.M., Urban, J., Rabe, F. (eds.) *Intelligent Computer Mathematics*. pp. 30–44. *Lecture Notes in Computer Science*, Springer, Berlin, Heidelberg (2011). https://doi.org/10.1007/978-3-642-22673-1_3
9. Huang, X., Kwiatkowska, M., Wang, S., Wu, M.: Safety verification of deep neural networks. In: *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part I. Lecture Notes in Computer Science*, vol. 10426, pp. 3–29 (2017)
10. Katz, G., Barrett, C., Dill, D., Ju-lian, K., Kochenderfer, M.: Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In: CAV (2017)
11. Katz, G., Huang, D.A., Ibeling, D., Julian, K., Lazarus, C., Lim, R., Shah, P., Thakoor, S., Wu, H., Zeljic, A., Dill, D.L., Kochenderfer, M.J., Barrett, C.W.: The Marabou framework for verification and analysis of deep neural networks. In: CAV 2019, Part I. LNCS, vol. 11561, pp. 443–452. Springer (2019)
12. Kokke, W., Komendantskaya, E., Kienitz, D., Atkey, R., Aspinall, D.: Neural networks, secure by construction - an exploration of refinement types. In: *Programming Languages and Systems - 18th Asian Symposium, APLAS 2020, Fukuoka, Japan, November 30 - December 2, 2020, Proceedings. Lecture Notes in Computer Science*, vol. 12470, pp. 67–85. Springer (2020)
13. Kokke, W., Komendantskaya, E., Kienitz, D., Atkey, R., Aspinall, D.: Neural Networks, Secure by Construction: An Exploration of Refinement Types. In: *Programming Languages and Systems*, vol. 12470, pp. 67–85. Springer International Publishing, Cham (2020). https://doi.org/10.1007/978-3-030-64437-6_4, series Title: *Lecture Notes in Computer Science*
14. Maria, E.D., Bahrami, A., L’Yvonnet, T., Felty, A.P., Gaffé, D., Ressouche, A., Grammont, F.: On the use of formal methods to model and verify neuronal archetypes. *Frontiers Comput. Sci.* **16**(3), 163404 (2022)
15. Passmore, G.O., Cruanes, S., Ignatovich, D., Aitken, D., Bray, M., Kagan, E., Kanishev, K., Maclean, E., Mometto, N.: The imandra automated reasoning system (system description). In: *Automated Reasoning - 10th International Joint Conference, IJCAR 2020, Paris, France, July 1-4, 2020, Proceedings, Part II*. vol. 12167, pp. 464–471. Springer (2020)
16. Passmore, G.O.: Some lessons learned in the industrialization of formal methods for financial algorithms. In: *Formal Methods - 24th International Symposium, FM 2021, Virtual Event, November 20-26, 2021, Proceedings. Lecture Notes in Computer Science*, vol. 13047, pp. 717–721. Springer (2021)
17. Sill, J.: *Monotonic Networks*. California Institute of Technology (1998)

18. Singh, G., Gehr, T., Püschel, M., Vechev, M.T.: An abstract domain for certifying neural networks. *PACMPL* **3**(POPL), 41:1–41:30 (2019). <https://doi.org/10.1145/3290354>
19. Wehenkel, A., Louppe, G.: Unconstrained monotonic neural networks. In: *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. pp. 1543–1553 (2019)
20. Wood, J.: Vectors and Matrices in Agda (Aug 2019), <https://personal.cis.strath.ac.uk/james.wood.100/blog/html/VecMat.html>