

# Learning with Boolean threshold functions

Veit Elser\*

*Department of Physics, Cornell*

Manish Krishan Lal†

*Department of Mathematics, Technische Universität München*

(Dated: February 20, 2026)

We develop a method for training neural networks on Boolean data in which the values at all nodes are strictly  $\pm 1$ , and the resulting models are typically equivalent to networks whose nonzero weights are also  $\pm 1$ . The method replaces loss minimization with a nonconvex constraint formulation. Each node implements a Boolean threshold function (BTF), and training is expressed through a divide-and-concur decomposition into two complementary constraints: one enforces local BTF consistency between inputs, weights, and output; the other imposes architectural concurrence, equating neuron outputs with downstream inputs and enforcing weight equality across training-data instantiations of the network. The reflect–reflect–relax (RRR) projection algorithm is used to reconcile these constraints.

Each BTF constraint includes a lower bound on the margin. When this bound is sufficiently large, the learned representations are provably sparse and equivalent to networks composed of simple logical gates with  $\pm 1$  weights. Across a range of tasks—including multiplier-circuit discovery, binary autoencoding, logic-network inference, and cellular automata learning—the method achieves exact solutions or strong generalization in regimes where standard gradient-based methods struggle. These results demonstrate that projection-based constraint satisfaction provides a viable and conceptually distinct foundation for learning in discrete neural systems, with implications for interpretability and efficient inference.

## I. INTRODUCTION

This year marks the 40th anniversary of “Learning representations by back-propagating errors”[1], the paper that launched the modern era of machine learning, known to the public as AI. Rumelhart, Hinton, and Williams understood that nonlinearities were required to break out of the more limited representations used by “perceptrons”, and that gradient optimization of network parameters could still be carried out efficiently, thanks to the chain rule of calculus. Though technical in nature, most workers in machine learning would find it hard to imagine a science of automated learning that was not built on the foundation of back-propagation.

Just three years prior to back-propagation, an optimization breakthrough in a little known image processing task called *phase retrieval* was making waves. Like artificial neural networks, phase retrieval works with large sets of data in the presence of nonconvexity. In neural networks, it is the function being minimized that is nonconvex, while in phase retrieval, it is the geometry of the constraint sets. Though constraint-based methods—with projections replacing gradient steps—had been used in phase retrieval for years, Jim Fienup’s “hybrid-input-output” algorithm [2] was miraculous by comparison. Whether neural network optimization or phase retrieval is the harder instance of nonconvexity is open to debate. In any case, it is interesting that the function minimiza-

tion approach never gained traction in phase retrieval (with constraints expressed in terms of an error function).

The constraint formulation of neural networks is no less natural than the usual one based on loss functions. In this article, we explore the alternate universe, where Fienup’s breakthrough in constraint satisfaction for nonconvex problems was so widely known that it was normal to consider it as an alternative to gradient descent—even by cognitive scientists. Because there is considerably more freedom in setting up constraint formulations than writing down loss functions, much of how we understand the new technique was developed for diverse applications—puzzles, sphere packings—that have nothing to do with phase retrieval or neural networks. This article contains only a brief summary. For a more comprehensive review suitable for nonexperts we refer the reader to [3].

This article is organized as follows. In the overview, we give a high-level comparison of the two approaches to network optimization, highlighting the advantages brought by the constraint-based approach. A key technical element for achieving interpretability and boosting generalization is the choice of nonlinear activation: the Boolean threshold function. In the review of BTFs that follows, it becomes clear that these are not “just functions” in the new setting. How the many constraints in neural network training are combined into just two easy constraints  $A$  and  $B$ , using the divide-and-concur framework, is covered next. This is followed by a brief exposition of Fienup’s algorithm and its generalization in geometric terms, where the work in each step (iteration) is performed by projections to the sets  $A$  and  $B$ .

---

\* ve10@cornell.edu

† manish.krishanlal@tum.de

A series of six applications follows. The problem of inferring a small multiplier circuit from a multiplication table is a good way to get familiar with the new approach. A historically interesting application comes next, where we revisit the binary autoencoder and manage to get true binary codes. MNIST-related problems follow, because the introduction of a new method would be incomplete without them. The applications that deserve the most attention are the ones that demonstrate generalization: logic circuits and cellular automata. The level of generalization is so striking for these synthetic data sets we know of nothing comparable in the function-minimizing approach to learning.

Readers eager to get a more hands-on experience can go to the software repository BOOLEARN and perform all the experiments reported here. The software is entirely in C, uses only the standard libraries, and runs in the Linux environment.

## II. OVERVIEW

The function-minimizing approach to neural network optimization places restrictions on admissible activation functions. In order to compute gradients, these must at least be differentiable almost everywhere. The sigmoid and hyperbolic-tangent functions not only are smooth everywhere, but they additionally have the nice feature that they saturate to values that can be interpreted as TRUE/FALSE.

The binary character of sigmoids was mentioned in the technical report [4] on which “Learning representations by ...” was based. Rumelhart et al. used an autoencoder network to drive home their optimization method’s ability to find weights that could map  $2^k$  linearly independent vectors to themselves, represented on  $2^k$  nodes, even when forced to pass through a “code-layer” with only  $k$  nodes. Nonlinear activation is essential for this parlor trick. The authors went on to remark that the 0/1 saturations of their sigmoid functions meant their network was encoding into and decoding-from, a binary code. But as their results showed, their codes often included the value  $1/2$ —the output of a sigmoid when it is furthest from being saturated.

The augmentation of the binary code can be overlooked when the application does not call on that level of interpretability. On the other hand, nice things are possible when the binary code is strict. The decoder-half of the autoencoder is a simple generative model, where  $2^k$  data vectors are generated by sampling a distribution on the  $k$  inputs. But the generative model only works when the admissible input-codes are a known code!

The easiest way to realize strict binary activation is through the Boolean threshold function:

$$y = \text{sgn}(w \cdot x). \quad (1)$$

This is the function used throughout this work. The weights  $w$  are the parameters of the network, and are

learned as in standard neural networks. But the output  $y$  and the elements of the input vector  $x$  are always  $\pm 1$ .

One way that Boolean threshold functions, or BTFs, are used differently here is that in addition to constraint (1), we also impose

$$|w \cdot x| \geq \mu. \quad (2)$$

The positive number  $\mu$  is the BTF’s *margin*. Our BTFs are not functions in the usual sense. The margin constraint (2), when satisfied with suitable weights for all the training data, compels all the inner products  $w \cdot x$  to avoid a gray area set by  $\mu$ . A successfully trained BTF network makes decisions with conviction.

Implementing BTFs as just described, with a loss function, while not impossible, is complicated and requires additional parameters. In our approach, constraints (1) and (2) are treated as a single “BTF constraint” and there is an elementary operation, called a projection, that takes an arbitrary trio  $(w, x, y) \in \mathbb{R}^{m+m+1}$  and outputs the unique trio  $(w', x', y')$  that satisfies the BTF constraint while minimizing

$$\|w' - w\|^2 + \|x' - x\|^2 + (y' - y)^2. \quad (3)$$

These projections are carried out concurrently on all the neurons/BTFs of the network. We call the independent, aggregate constraints at the BTFs, the *A* constraint. Geometrically, *A* is a continuous set in the high-dimensional space of  $w$ ,  $x$  and  $y$  variables. When the network has  $N$  nodes and  $E$  edges, the dimension of that space is  $N+2E$ . The weights ( $w$ ) and inputs ( $x$ ) live on the edges because they are associated with sending/receiving node-pairs.

The idea behind treating all the neurons independently, in the *A* constraint, is called *divide-and-concur* [5]. There is another aggregate constraint, called *B*, that imposes “concur” or the equality of one neuron’s output ( $y$ ) with another neuron’s input ( $x$ ). Implementing the network architecture is strictly the job of the *B* constraint. This constraint is a hyperplane (a system of many linear equations) and is even easier to project to than the *A* constraint.

The weights also appear in the *B* constraint, in a way that represents another departure from gradient-based optimization. Instead of there being one set of network variables (comprising  $w$ ’s,  $x$ ’s, and  $y$ ’s), the constraint-approach uses as many instantiations of the network as there are data items. Each instantiation has its own set of  $x$  and  $y$  variables, since each data item constrains these differently at the network’s inputs and outputs. But forcing the  $w$  variables to be the same across all the instantiations would greatly complicate the projection to the *A* constraint. Not only do we want  $(w, x, y)$  to be independent across all the nodes/BTFs, we want that independence to extend across data items. The divide-and-concur solution is to have separate weights for each data item (to keep the *A* constraint simple) and impose their equality in the *B* constraint.

In addition to forcing the weights to concur across data items, the  $B$  constraint imposes the constraint

$$w \cdot w = m \quad (4)$$

on a BTF's weights when it has  $m$  inputs. This matches the normalization of the Boolean inputs,  $x \cdot x = m$ . And since  $y = \pm 1$ , all three variable types have rms value 1. Sensible variable normalizations are important if we use the distance (3) when computing projections.

Gradient-descent optimized networks are trained in the “stochastic” mode, where the gradient of the loss is evaluated for small, random batches of the data. When the entire data set is used, there is too much cancellation in the gradient computation for the optimizer to act on. There is no analogous problem in the constraint-based approach.

To understand the last statement, one needs to first learn how the constraint satisfaction algorithm, using projections, finds variable assignments that satisfy both  $A$  and  $B$ . This is covered in Section V. For the present, it suffices to know that if  $z_A$  and  $z_B$  are projections of the full set of network variables—denoted by  $z$ —to the two constraints, then the optimizer moves  $z$  by the rule

$$z' = z + \beta(z_A - z_B)$$

in each iteration. The positive number  $\beta$  is the “time-step” and is roughly analogous to the learning-rate parameter of gradient descent. But the difference  $z_A - z_B$  is not the gradient of anything. Instead, the “gap”

$$\Delta = \|z_A - z_B\| \quad (5)$$

measures the distance that is left to reconcile between the “divide” and “concur” constraints.

When the gap is zero, both constraints are satisfied. That is, all the individual BTF constraints are satisfied, for all the data items, and there is just one shared set of weights. In constraint problems where  $A$  and  $B$  are convex,  $\Delta$  always decreases and behaves like the loss in gradient descent [6]. But in our case  $A$  is not convex, and  $\Delta$  can increase. When this happens, the algorithm is extricating itself from situations where the  $A$  and  $B$  variable assignments are locally irreconcilable. This nice feature persists when the training batches are large. So unlike stochastic-gradient-descent, which requires small batches, in the constraint-based method the batch size should be as large as the computing resources allow.

The margin (2) and weight normalization (4) constraints, acting together, bring about another nice feature: network sparsity. Like the strict Boolean activation of the BTF, networks with few relevant edges (having significant weight) bring interpretability to the representation. In Section III we prove that by defining the imposed margin on  $m$ -input BTFs as

$$\mu_m = \sqrt{m/\sigma}, \quad (6)$$

where  $\sigma$  is called the *support hyperparameter*, then up to a caveat all the BTFs in the network will be equivalent to

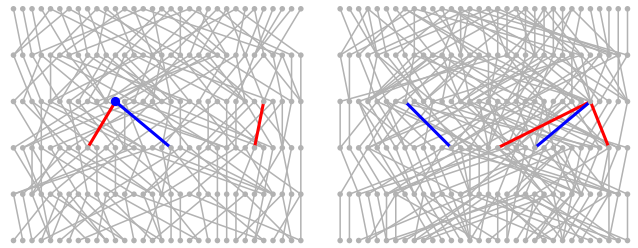


FIG. 1. Two 5-layer random logic circuits where each node either copies a value from the layer below (isolated colored lines) or applies a logical gate to them (colored groups). The gates are either 2-input AND/OR (left circuit) or 3-input MAJ (right circuit). The colors (red/blue) of the highlighted edges indicate whether or not negation is applied. The blue node in the left circuit means this BTF is taking input from the constant node with weight +1, thereby choosing to implement AND.

BTFs having at most  $\sigma$  nonzero elements in their weight vectors. The caveat is that, in training, each BTF must see the complete set of possible input patterns on its “relevant” weights—those that are nonzero in the equivalent BTF. This condition has a reasonable chance of being satisfied when the BTF supports are small.

A corollary of the support property that we impose through the margin constraint and (6), is the property that when  $\sigma$  is an odd integer the only way that BTFs can satisfy the equality case of the margin constraint is when the weight vector has exactly  $\sigma$  nonzero and equal-magnitude elements. The case  $\sigma = 3$ , corresponding to weights

$$w = \sqrt{m/3}(\pm 1, \pm 1, \pm 1, 0, \dots, 0) \quad (7)$$

and permutations, is already interesting. BTFs with such weights implement the 3-input majority gate (MAJ), which acts as a 2-input AND or OR when one of its inputs is held fixed. All the BTFs in our networks have an edge to a special node with value  $-1$ , so by choice of the weight to that node, the BTF can elect to implement a 2-input AND/OR. The margin  $\mu_m = \sqrt{m/3}$  also admits BTFs with fewer than 3 relevant input weights, but these are all equivalent to a BTF with just one relevant input—the COPY-gate.

There is no better demonstration of what the constraint-based method is capable of than the reconstruction of logic circuits (Fig. 1). Whereas the number of Boolean functions from  $m$  inputs to  $m$  outputs is astronomical ( $2^{m2^m}$ ), learning the function from a modest amount data is possible, in principle, if we know it can be expressed by a circuit comprising just NOT and few-input logic gates arranged in a few-layer network. When deployed on a network with fully connected layers and the setting  $\sigma = 3$ , our method is forced to consider only circuits where each BTF either just copies a Boolean value from one of the nodes in the layer below, or applies simple logic to 2 or 3 of those values (depending on

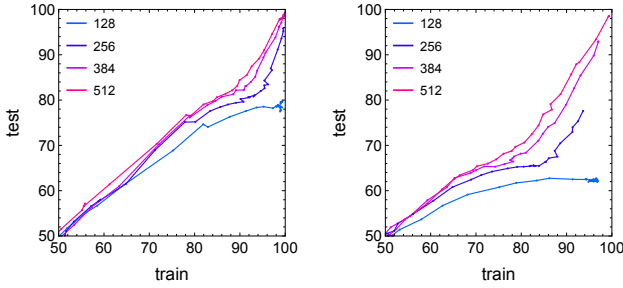


FIG. 2. Evolution of the training and test accuracies when training on 128, 256, 384, and 512 data generated by the random logic networks above. Learning the random AND/OR data (left plot) is easier, but 256 data appear to be sufficient to learn both types of data, given sufficient iterations of the algorithm.

whether the special constant node was chosen as one of the three allowed by  $\sigma = 3$ ). In any case, there is no information telling the constraint solver how one layer is sparsely wired to the next, or where to insert the NOT gates (weights of negative sign).

While more details are given in Section VI E, a preview of circuit reconstruction should convey that constraint-based learning with BTF networks is remarkable. Figure 1 is a rendering of two 5-layer circuits that were used to generate data—pairs of strings of 32 bits. The evolution of the training and test accuracies for different numbers of training data are plotted in Figure 2. Accuracies are just the fraction of correct output bits when given an input string from the training set or, for testing, a string the learning algorithm has never seen. We see evidence of a transition at 256 training items. Below this number the reconstructed circuits are clearly different from those in Figure 1 because the test accuracy is well below 100%. But above 256 training items the test accuracies for both data sets appear to be on track to reach 100%, given enough iterations of the constraint solver (limited to  $10^6$  in this demonstration).

Figure 3 shows accuracy results for the same data when the gradient-descent method is used. All the results we report using this method should be seen as *baselines*: starting points that use state-of-the-art optimization on the most widely used multilayer perceptron model. Details can be found in appendix A. Refinements of the baseline method clearly have a ways to go to equal the performance of the constraint-based method. The test accuracies in Figure 3 show no signs of reaching 100%, and there is no evidence of a transition at some number of training data. We estimate the information sufficiency for perfect circuit reconstruction in Section VI E and find it is just above 128 items. While both methods were limited to  $10^6$  steps/iterations, the plots suggest only the constraint method’s test accuracy would improve when this number is increased.

High test-accuracies, when faced with “a poverty of stimulus,” can be explained by the fact that the learning

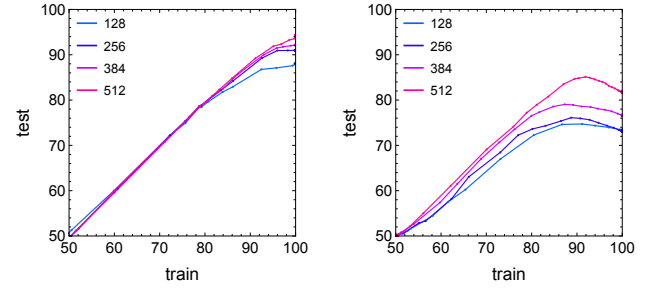


FIG. 3. Same evolution of learning curves as in Figure 2 but for the gradient-descent method.

algorithm is forced to only consider models that are very close to the models that generated the data (same number of layers, sparse, etc.). That would be truly amazing and the use of the term “reconstruction” would be justified. However, the distribution of the learned weights shows that something not quite that simple is going on. These are shown for the 5-layer-circuit data in Figure 4. Because the weight vector for every neuron has normalization  $w \cdot w = 32 + 1$ , when we see peaks in the distribution at  $\pm\sqrt{33/1} \approx \pm 5.7$  it means there are neurons receiving effectively just one input (COPY). Similarly, peaks at  $\pm\sqrt{33/3} \approx \pm 3.3$  correspond to AND/OR/MAJ gates. The distributions in the lower panels show just the lowest-layer weights, and these are indeed close to the values just described. The upper panels show the weight distributions over all five layers, and these are more diverse. Is the algorithm getting sloppy in the higher layers?

The diversification of the weights has an innocent explanation. In the ground truth models (Figure 1) there is a small probability that the values at two nodes in the same layer are perfectly correlated, or perfectly anticorrelated (when these are copied from the same node in the layer below). A gate taking input from one of these nodes could instead take input from the other node (negated, if anticorrelated) without changing its output. But it could also take a fixed mixture, with weights  $pw$  and  $(1-p)w$ , and still produce the same output. This has the effect

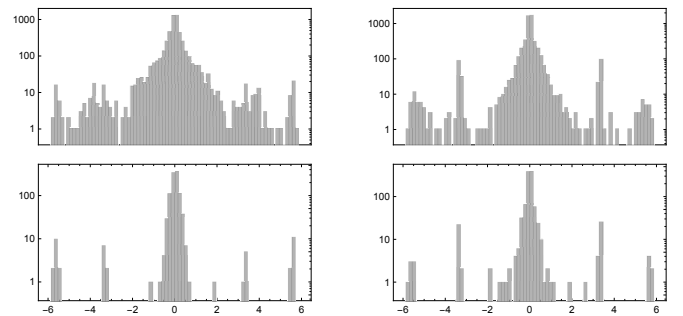


FIG. 4. Distribution of network weights after training on the two types of random logic data (left: AND/OR, right: MAJ). The lower histograms show the weights in just the first layer of each network.



of decreasing the neuron's  $w \cdot w$  norm, with the even mixture ( $p = 1/2$ ) giving the greatest decrease. To restore the norm, the constraint solver will scale up all of the neuron's weights, thereby also increasing the margin. The consolidation of nodes, as in this example of three relevant nodes becoming four, is inevitable when there is a loss of entropy in the higher layers of the network. It allows for a diversification of the network weights without compromising functionality.

Learning a Boolean function from examples is known as *Boolean function inference*. In machine learning the same task would be described as generalization, albeit with rather strong constraints on the model. Models that implement simple logic operations over some small number of layers might be a good fit for various kinds of symbolic data. This article is not about those potential applications, but about a technology that can train such models.

### III. BOOLEAN THRESHOLD FUNCTIONS AS CONSTRAINTS

A Boolean threshold function, or BTF, is a function  $f : \{-1, 1\}^m \rightarrow \{-1, 1\}$  that can be represented by a weight vector  $w \in \mathbb{R}^m$  as

$$f(x) = \text{sgn}(w \cdot x). \quad (8)$$

We will use  $X_m = \{-1, 1\}^m$  for the hypercube of possible inputs and restrict weights as “admissible” by the property

$$\forall x \in X_m : w \cdot x \neq 0. \quad (9)$$

In the machine learning context we will refer to the values  $\pm 1$  as “Boolean”, where  $-1$  corresponds to the more standard 0 or FALSE. In networks, a weight vector's negative elements act as NOT gates.

BTFs have been the subject of much research, mostly in connection with the *Chow parameter problem* [7]. This is the problem of constructing a representation—an admissible  $w$ —when given the set  $X_+ = f^{-1}(1)$  for some BTF (more specifically, parameters derived from  $X_+$  first introduced by Chow). Thankfully, this difficult chapter of BTF history is largely irrelevant for the margin and support characteristics of BTFs that this work relies on.

Since we are always interested in the representations, we use the notation  $f_w$  for BTFs when a question involves the representation/weights  $w$ . The number of inputs is always  $m$ . The property of BTFs that we will exploit in networks calls on the following definitions.

**Definition III.1.** *The support of a weight vector  $w$ , or  $\text{supp}(w)$ , is the number of “relevant” inputs of  $f_w$ . Input  $p$  is relevant if and only if there exists some  $x \in X_m$  such that*

$$f_w(x_1, \dots, x_p, \dots, x_m) \neq f_w(x_1, \dots, -x_p, \dots, x_m).$$

**Definition III.2.** *The margin of a weight vector  $w$  is given by*

$$\mu(w) = \min_{x \in X_m} |w \cdot x| \quad (10)$$

*and is strictly positive when  $w$  is admissible.*

A BTF with  $\text{supp}(w) = 1$  just acts as a wire that conveys a single input value to the output, with or without negation. There are no BTFs with  $\text{supp}(w) = 2$ . That's because if  $w = (w_1, w_2)$ , then to be admissible  $|w_1| \neq |w_2|$ . But the output of such a BTF only depends on the input having the largest magnitude weight, so  $\text{supp}(w) = 1$ . BTFs first become interesting at  $\text{supp}(w) = 3$ .

Fully supported weight vectors have the following bound on their margin:

**Lemma III.3.** *If  $w$  is admissible and  $\text{supp}(w) = m$ , then*

$$\mu(w) \leq w_{\min} = \min_p |w_p|. \quad (11)$$

*Proof.* Let  $p$  be an input for which  $|w_p|$  is a minimum and call this number  $w_{\min}$ , where  $w_{\min} > 0$  since otherwise  $\text{supp}(w) < m$ . Now consider the set

$$Y_{-p} = \{w \cdot x - w_p x_p : x \in X_m\}, \quad (12)$$

that is, the set of threshold function values without the contribution from input  $p$ . At least one element  $y^* \in Y_{-p}$  satisfies  $|y^*| < w_{\min}$ , since otherwise there are no elements  $y \in Y_{-p}$  for which  $y \pm w_{\min}$  have opposite signs and input  $p$  would be irrelevant (so  $w$  would not have full support). But  $|y^* \pm w_{\min}|$  are bounds on the margin, so we find

$$\mu(w) \leq |y^*| - w_{\min} \leq w_{\min}. \quad (13)$$

□

Because the weight vectors in our networks have a fixed normalization, the result we rely on the most is the following:

**Theorem III.4.** *Let  $w$  be an admissible weight vector with  $\text{supp}(w) = m$ , then*

$$\frac{\mu(w)}{\|w\|} \leq \frac{1}{\sqrt{m}} \quad (14)$$

*and the equality case requires that  $m$  is odd and all the elements of  $w$  have equal magnitude.*

*Proof.* Using lemma III.3 on the numerator,

$$\frac{\mu(w)}{\|w\|} \leq \frac{w_{\min}}{\sqrt{\sum_{p=1}^m w_p^2}} \quad (15)$$

$$\leq \frac{w_{\min}}{\sqrt{mw_{\min}^2}} = \frac{1}{\sqrt{m}}. \quad (16)$$

Equality requires that all  $m$  terms in the denominator are equal, or that the weight vector has equal-magnitude elements. But such a weight vector for even  $m$  is not admissible, so the equality case is only possible for odd  $m$ .  $\square$

When weight vectors  $w$  and  $w'$  represent the same BTF we write  $w \sim w'$ . In our networks the number of BTF inputs  $m$  may vary, and the weight vectors are normalized as

$$\|w\|^2 = m. \quad (17)$$

With this choice, the rms value of the weight vector elements matches that of the Boolean inputs  $x$  and the Boolean output  $y$ . Under the normalization constraint, maximizing the margin translates to fewer nonzero elements:

**Lemma III.5.** *Suppose an  $m$ -input BTF has  $\text{supp}(w) = n \leq m$ . Then there exists a  $w^* \sim w$  with exactly  $n$  nonzero weights and  $\mu(w^*) \geq \mu(w)$ .*

*Proof.* The weights of  $w$  associated with the  $m - n$  irrelevant inputs can be set to zero without changing the BTF. To restore normalization (17) to the resulting  $w^*$ , the nonzero elements are rescaled by a factor  $\lambda \geq 1$ . This changes the margin by the same factor.  $\square$

This property of BTF representations combined with theorem III.4 will let us learn network representations of Boolean data that are sparse.

### A. Constraints on margins in training

During training each  $m$ -input BTF with weight vector  $w$  is subject to the following constraint:

$$\forall x \in X_t : \quad \mu_m \leq |w \cdot x|. \quad (18)$$

Here  $\mu_m$  is an  $m$ -dependent margin lower bound set by the user, and  $X_t \subset X_m$  is the set of inputs seen by the BTF during training.  $X_t$  might be a proper subset because the size of the training data is small, the node consolidation effect, or a combination of these. Because of this possibility it may be possible to satisfy all the margin constraints (18) even when  $\mu_m > \mu(w)$ .

The possibility that  $X_t$  is significantly smaller than  $X_m$  is made unlikely by a combination of two things. First, in constraint-based training one is not forced to use small data batches as in the stochastic approach to gradient-based training. The second reason is slightly circular. When sparse network representations exist, where the BTFs all have small support, then it is only the size of the subset of the relevant inputs that matter. A BTF with support  $n$  much smaller than  $m$  only has to see the much smaller complete set of  $2^n$  vectors on the relevant inputs in order to satisfy the more generous constraint

$$\mu_m \leq \mu(w^*) \quad (19)$$

promised by lemma III.5 for the representation  $w^*$  that has only  $n$  nonzero elements.

For the reasons just given, it makes sense to define the  $m$ -dependent margin bound as

$$\mu_m = \sqrt{m/\sigma}, \quad (20)$$

where  $\sigma$  is the *support hyperparameter*. Assuming a sparse network representation exists, with sparse BTF representations  $w^*$ , then applying theorem III.4 to a BTF with  $n \leq m$  nonzero elements in its weight vector  $w^*$ , we find

$$\mu(w^*) \leq \frac{\|w\|}{\sqrt{n}} = \sqrt{m/n}. \quad (21)$$

Theorem III.4 only applies when the inputs to the BTF are sufficiently diverse as described above, but if that is true, then (19) is also true and we arrive at

$$n \leq \sigma. \quad (22)$$

We now have a strategy for learning sparse network representations: force the BTF margins to be large with a small  $\sigma$ .

In our applications of the new training method we will encounter situations where the bound (22) is rigorous. When the network has  $m$  inputs and data for all  $2^m$  unique inputs to the first layer of BTFs is available, then in any solution to the constraint satisfaction problem those BTFs will have supports  $n$  that are upper bounded by the parameter  $\sigma$ . In the event that the smallest  $\sigma$  for which solutions exist is an odd integer  $n^*$ , then theorem III.4 also tells us that for the BTFs (in the first layer) with exactly  $n^*$  nonzero weights, these will have equal magnitudes.

### B. Representing logic gates

The weights  $w$  of BTFs have simple structures when their margin  $\mu(w)$  is maximized subject to a fixed norm constraint. The maximizing  $w^*$  is uniquely determined by no more than  $m$  “extreme” inputs  $x$  of  $f_w^{-1}(1)$  that have a common value of  $w^* \cdot x$ , and this value equals  $\mu(w^*)$ . Not surprisingly, the margin-maximizing  $w^*$  are no different from the linear programming optimized weights that were tabulated as early as 1961 [8]. Up to support  $m = 8$  these are integer vectors when normalized to have  $\mu(w) = 1$ . For  $m = 9$  there are counterexamples, such as

$$w = \frac{1}{2}(29, 25, 19, 15, 12, 8, 8, 3, 3). \quad (23)$$

To admit a BTF with this  $w$ , our “support” hyperparameter would have to satisfy  $\sigma \geq \|w\|^2 = 1171/2$ .

Packing highly complex decision encodings into individual neurons was not our motivation for using BTFs! Not only would the small margin constraints  $\mu_m = \sqrt{m/\sigma}$  be problematic for robust learning, they would

admit BTFs whose actual supports are very large. By keeping  $\sigma$  small, the complexity of the representation is transferred from the processing internal to neurons, to the complexity of the wiring between them.

To add perspective to the objective just described, and defend our choice of the term “support”, consider the case  $\sigma = 9$ . The type of a BTF can be uniquely specified by a nonincreasing tuple of positive weights as in (23). The full set of BTFs of a given type is obtained by applying all possible sign changes and permutations to the positive and ordered tuple. Here is the complete inventory of BTF types that are admitted with  $\sigma = 9$  (squared-norm 9 or less) :

$$\begin{aligned} &(1, 0, \dots) \\ &(1, 1, 1, 0, \dots) \\ &(1, 1, 1, 1, 1, 0, \dots) \\ &(2, 1, 1, 1, 0, \dots) \\ &(1, 1, 1, 1, 1, 1, 0, \dots) \\ &(1, 1, 1, 1, 1, 1, 1, 0, \dots) . \end{aligned}$$

In all cases except one, the value of  $\sigma$  translates directly to the support of the BTF, and the BTF implements a simple majority gate. We should not be dismissive of the simplicity of majority gates. After all, the 9-input gate can take 512 forms depending on the negation pattern of its weights.

In most of our applications we will set  $\sigma = 3$ . This is sufficient to implement arbitrary Boolean functions in a deep network because the 3-input majority gate can implement 2-input AND/OR:

$$\begin{aligned} \text{MAJ}(-1, x_1, x_2) &= \text{AND}(x_1, x_2) \\ \text{MAJ}(+1, x_1, x_2) &= \text{OR}(x_1, x_2) . \end{aligned}$$

Our definition of BTFs with threshold value zero makes them *self-dual*. Self-dual Boolean functions  $f$  have the property  $f(x) = -f(-x)$  for all inputs  $x$ . AND/OR (for any number of inputs) are not self-dual because they are expressed in terms of a self-dual function with one input held constant. The BTFs in our networks all take input from a special node 0 with value  $x_0 = -1$ . By choice of the corresponding weight  $w_0$ , a BTF can choose to be self-dual ( $w_0 = 0$ ), and-like ( $w_0 > 0$ ), or or-like ( $w_0 < 0$ ). The special weight  $w_0$  plays the same role as the bias parameter in standard neural networks.

Keeping a single input constant does not change any of the results in this section that used a margin defined in terms of a minimum over the entire hypercube of inputs. That’s because

$$\begin{aligned} \min_{x \in X_m} |w \cdot x| &= \min_{x \in X_m^- \cup X_m^+} |w \cdot x| \\ &= \min \left( \min_{x \in X_m^-} |w \cdot x| , \min_{x \in X_m^+} |w \cdot x| \right) \\ &= \min \left( \min_{x \in X_m^-} |w \cdot x| , \min_{x \in X_m^-} |-w \cdot x| \right) \\ &= \min_{x \in X_m^-} |w \cdot x| , \end{aligned}$$

where  $X_m^\pm$  is the half-hypercube whose points have  $\pm 1$  for their first coordinate.

### C. Projecting to the BTF constraint

The set  $A$  corresponding to the BTF constraint is the union  $A = A_+ \cup A_-$ , where (all upper or all lower signs)

$$A_\pm = \left\{ (w, x, y) \in \mathbb{R}^{m+m+1} : \pm w \cdot x \geq \mu_m, y = \pm 1 \right\} .$$

Projecting to the BTF constraint means, for any given  $(w, x, y) \in \mathbb{R}^{m+m+1}$ , compute the point  $(w', x', y') \in A$  that minimizes the distance (3). The projection is unique except for  $(w, x, y)$  on a set of measure zero.

There are four cases, depending on whether  $w \cdot x$  and  $y$  have the same sign, and whether  $|w \cdot x|$  exceeds the margin. For the easiest case, where the signs match and  $|w \cdot x| \geq \mu_m$ , the projection only needs to change  $y : y' = \text{sgn}(y)$ . In the most compute-intensive case, where both of these conditions are violated, two outcomes must be considered:

$$\begin{aligned} w' \cdot x' &= +\mu_m, & y' &= +1 \\ w' \cdot x' &= -\mu_m, & y' &= -1 . \end{aligned}$$

The projection of  $(w, x)$  to these bilinear constraints can be worked out using the method of Lagrange multipliers. The result,

$$\begin{aligned} w' &= \frac{w + \lambda x}{1 - \lambda^2} \\ x' &= \frac{x + \lambda w}{1 - \lambda^2} , \end{aligned}$$

is expressed in terms of the unique zero  $\lambda \in (-1, 1)$  of the function

$$h_\pm(\lambda) = \frac{(1 + \lambda^2)w \cdot x + \lambda(w \cdot w + x \cdot x)}{(1 - \lambda^2)^2} \mp \mu_m ,$$

where the upper sign corresponds to the case  $y' = +1$ . In large networks the work to compute this projection scales as  $m$ , the number of vector elements that are updated and terms in the three inner product computations. An efficient root-finding method is described in [9, 10].

### D. Symmetries

In addition to any permutation symmetries of the network, BTF networks have another symmetry not shared by the networks of standard machine learning. This is a local symmetry at each of the network’s “hidden” nodes, defined to be nodes that both receive and send Boolean values to other nodes. This symmetry is a generalization of De Morgan’s law and follows directly from the self-dual property of the BTF:

$$\begin{aligned} f_w(x) &= -f_w(-x) \\ &= -f_{-w}(x) . \end{aligned}$$

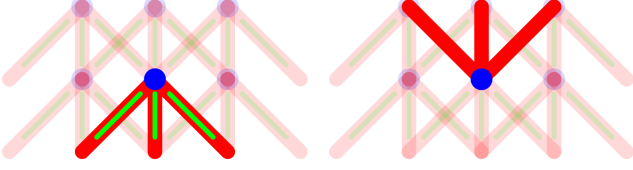


FIG. 5. Divide-and-concur variable groupings in BTF networks. Red lines are neuron inputs ( $x$ ), blue disks are neuron outputs ( $y$ ), and green lines are weights ( $w$ ). On the left, the highlighted fan-in shows all the variables in one BTF constraint. On the right, the highlighted fan-out shows the concur between neuron outputs and inputs.

Simultaneously flipping the signs of the weights on all the incoming and outgoing edges of a hidden node has no effect on the Boolean function the network is expressing.

#### IV. DIVIDE-AND-CONCUR FOR BTF NETWORKS

In this section we consider general feed-forward network architectures. There is no reference to layers, or the special node that holds the constant input.  $\mathcal{N}$  is the set of all nodes; hidden, input and output nodes are denoted  $\mathcal{H}$ ,  $\mathcal{I}$ , and  $\mathcal{O}$ . The network edges are denoted  $\mathcal{E}$  and  $\mathcal{D}$  is the set of data labels.

Variables have both a network and data label. At each node  $p$ , and for each data item  $i$ , there is an “output variable”  $y_{pi}$ . These hold network inputs when  $p \in \mathcal{I}$  and BTF outputs otherwise. The other variables live on network edges, labeled  $p \rightarrow q$  when the BTF at node  $q$  receives input from node  $p$ . The “input variables”  $x_{p \rightarrow qi}$  have this structure to allow them to satisfy the constraints for BTF  $q$  without having to be consistent with  $y_{pi}$ . The “weight variables”  $w_{p \rightarrow qi}$  clearly also have an edge label, and  $i$  labels the data-item copy that allows the BTF constraints to be satisfied independently for each  $i$ .

The groupings of the three variable types in the divide-and-concur framework are shown in Figure 5. Highlighting in the left panel shows “divide”: a disjoint union by BTFs, each one a fan-in of inputs (red lines), weights (green lines), with a single BTF output (blue disk). The constraints on each group are independent and together comprise the  $A$  constraint. The highlighted fan-out in the right panel shows a different disjoint union. The constraint here is the concurrence of one BTF output with all the inputs that receive that output. Taken together, the independently concurring groups are part of the  $B$  constraint. The  $B$  constraints also includes the concurrence of weights across data items (green lines across multiple network instantiations).

More formally,  $A$  has the following product structure

$$A = \prod_{q \in \mathcal{H} \cup \mathcal{O}}^\times \prod_{i \in \mathcal{D}}^\times A_{qi}.$$

Each factor has the form given in Section III C with  $y$  replaced by  $y_{qi}$ ,  $w$  and  $x$  by the vectors with elements  $w_{p \rightarrow qi}$  and  $x_{p \rightarrow qi}$ , where  $p$  ranges over all nodes incident on  $q$ , which we denote  $p \in \rightarrow q$ . The  $A$  constraint also has a factor associated with the input nodes. This is described later, as it involves the data.

The output  $y_{pi}$  of BTF  $p$  for data item  $i$  and the inputs  $x_{p \rightarrow qi}$  sent to BTFs  $q$  (for the same data item) are no longer independent in the “concur” or  $B$  constraint. Even so,  $B$  also has a product structure:

$$B = \prod_{p \in \mathcal{I} \cup \mathcal{H}}^\times \prod_{i \in \mathcal{D}}^\times B_{pi}.$$

Each  $B_{pi}$  is an independent constraint defined by the system of equalities

$$\forall q \in p \rightarrow : x_{p \rightarrow qi} = y_{pi}, \quad (24)$$

where  $p \rightarrow$  denotes the set of nodes on which  $p$  is incident. The  $B$  constraint also has a factor associated with the data (described later), now at all the output nodes.

To keep the weights across all the data items consistent, the  $B$  constraint also has the following factors, one at each BTF:

$$\prod_{q \in \mathcal{H} \cup \mathcal{O}}^\times B_q.$$

Each factor labeled  $q$  combines weight-equality across data items and weight normalization:

$$\forall p \in \rightarrow q : w_{p \rightarrow qi} = \bar{w}_{p \rightarrow q}, \quad (25)$$

$$\sum_{p \in \rightarrow q} \bar{w}_{p \rightarrow q}^2 = |\rightarrow q|. \quad (26)$$

The first statement simply states that every data-item copy  $i$  is equal to the same “concur” value, with symbol  $\bar{w}_{p \rightarrow q}$ . The squared-norm equals the number of inputs to the BTF,  $|\rightarrow q|$  (denoted  $m$  in Section III).

The  $B$  constraint is so simple that with essentially no extra work we can compute the projection to this set for a natural generalization of the metric:

$$\|z' - z\|_g^2 = \sum_{q \in \mathcal{N}, i \in \mathcal{D}} g_{qi} \|z'_{qi} - z_{qi}\|^2, \quad (27)$$

where the symbol  $z$  is being used for the aggregate of  $(w, x, y)$  variables,

$$\|z'_{qi} - z_{qi}\|^2 = (y'_{qi} - y_{qi})^2$$

for  $q \in \mathcal{I}$ , and otherwise

$$\begin{aligned} \|z'_{qi} - z_{qi}\|^2 = \sum_{p \in \rightarrow q} & \left( (w'_{p \rightarrow qi} - w_{p \rightarrow qi})^2 \right. \\ & + (x'_{p \rightarrow qi} - x_{p \rightarrow qi})^2 \\ & \left. + (y'_{qi} - y_{qi})^2 \right). \end{aligned}$$

is the same as expression (3) of Section II (so that the projection to the BTF constraint is unchanged). The



positive numbers  $g_{qi}$  are the *metric coefficients*. For the input nodes  $q$ , and all data items  $i$ , we set  $g_{qi} = 1$ . However, we will see that allowing the metric coefficients to be nonuniform elsewhere enhances the performance of the constraint satisfaction algorithm. A simple, automated heuristic for tuning the metric is described in Section V. In this section we treat the metric coefficients as fixed, given numbers.

Derivations of the concur values for (24), in the case of a nonuniform metric, can be found in [3]:

$$\bar{y}_{pi} = \frac{\sum_{q \in p \rightarrow} g_{qi} x_{p \rightarrow qi} + g_{pi} y_{pi}}{\sum_{q \in p \rightarrow} g_{qi} + g_{pi}}.$$

For the projection of the weights in the  $B$  constraint one first computes

$$\bar{w}_{p \rightarrow q} = \frac{\sum_{i \in \mathcal{D}} g_{qi} w_{p \rightarrow qi}}{\sum_{i \in \mathcal{D}} g_{qi}},$$

and rescales this average by the factor

$$\sqrt{\frac{|\rightarrow q|}{\sum_{p \in \rightarrow q} \bar{w}_{p \rightarrow q}^2}}$$

to satisfy the side constraint (26).

### A. Input and output constraints

The data value imposed at network input  $p \in \mathcal{I}$ , for data item  $i$ , is written  $y_{pi}^A$  because this constant takes the place of the projection to the BTF constraint—called  $A$  above—at the input nodes (where there is no BTF). Similarly, the data value imposed at network output  $q \in \mathcal{O}$ , for data item  $i$ , is written  $y_{qi}^B$  because this constant takes the place of the projection to the concur constraint—called  $B$  above—at the output nodes (from which there are no  $x$  variables to concur with). The corresponding data factors of  $A$  and  $B$  are simply

$$A_{\mathcal{I}} = \left\{ y \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{D}|} : y_{pi} = y_{pi}^A \right\} \quad (28a)$$

$$B_{\mathcal{O}} = \left\{ y \in \mathbb{R}^{|\mathcal{O}| \times |\mathcal{D}|} : y_{qi} = y_{qi}^B \right\}. \quad (28b)$$

In BOOLEARN the data values in files are always 0 or 1, and these are converted to  $-1$  and  $+1$ , respectively, to be the  $y^A$  values used by the constraint satisfaction algorithm. There is one exception, when the input data is analog and sufficiently close to Boolean, like the pixel values of MNIST. In this case, BOOLEARN assumes the input data lies in the interval  $[0, 1]$  and these are mapped linearly into the interval  $[-1, 1]$  to give  $y^A$ .

The training program in BOOLEARN can be run in a generative mode, where no input data is given. Two options for constraining the inputs in that case are available. The weakest is the constraint

$$A_{\mathcal{I}} = \left\{ y \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{D}|} : y_{pi} \in \{-1, 1\} \right\}, \quad (29)$$

which imposes Boolean values when there is no BTF to do that. Projecting to this constraint means rounding to  $-1$  or  $1$ , independently for each input  $p$  and data item  $i$ . The other option is the 1-hot constraint:

$$A_{\mathcal{I}} = \left\{ y : y_{pi} \in \{-1, 1\}, \sum_{p \in \mathcal{I}} y_{pi} = 2 - |\mathcal{I}| \right\} \quad (30)$$

where  $y \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{D}|}$ . To project to this constraint, independently for each  $i$ , the largest input  $p$  is singled out and replaced by  $1$ , all others by  $-1$ .

## V. CONSTRAINT SATISFACTION WITH RRR

By using the divide-and-concur strategy, in Section IV we showed how the BTF-network training problem can be cast in the form

$$\text{find some } z \in A \cap B, \quad (31)$$

where  $z$  is a point in the space of the aggregated  $(w, x, y)$  variables,  $A$  is the set in that space that defines the operation of BTFs, and  $B$  defines the network’s architecture. The training data provides additional constraints, and these are included with  $A$  at the inputs and  $B$  at the outputs. In this section we describe an algorithm for solving (31) called reflect-reflect-relax (RRR), a generalization of Fienup’s hybrid-input-output phase retrieval algorithm. RRR makes no reference to the “divide” character of  $A$  and the “concur” character of  $B$ . In fact, there are many applications of RRR where  $A$  and  $B$  do not have those interpretations at all.

In order to use RRR, the one property required of  $A$  and  $B$  is that the two projections,

$$P_A(z) = \arg \min_{z' \in A} \|z' - z\| \\ P_B(z) = \arg \min_{z' \in B} \|z' - z\|,$$

can be computed efficiently. What we showed in Section IV is that not only can these be computed efficiently for the BTF-network training problem, but that the computations distribute, thanks to the product structures of  $A$  and  $B$ . In the case of  $P_A$ , the variables assigned to each BTF in the network and each item in the data are independently projected to satisfy the sgn “activation function” with its lower bound on the margin.  $P_B$  is distributed over all the groups of BTF outputs  $y_{pi}$  and receiving-BTF inputs  $x_{p \rightarrow qi}$ , and also independently over the data  $i$ .  $P_B$  also concurs the weights (across data), independently for each edge of the network.

$P_A$  and  $P_B$  are computed once in each iteration of RRR, and the work is proportional both to the number of network edges and the number of data items. The work per RRR iteration is therefore comparable to the work performed in gradient descent over a single pass through the data.

For the history of RRR and the logic behind the iterations rule,

$$z \mapsto z' = z + \beta \left( P_B(P_A(z)) - P_A(z) \right), \quad (32)$$

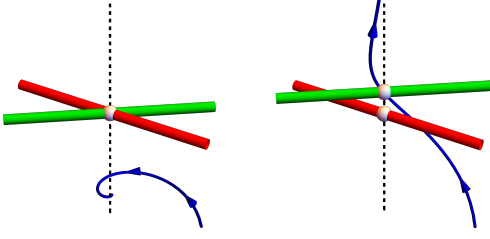


FIG. 6. Local behavior of the iterate  $z$  (blue) in the limit of small  $\beta$  when the problem is locally feasible (left) and locally infeasible (right).  $A$  and  $B$  can locally be approximated as affine and are rendered as the red ( $A$ ) and green ( $B$ ) rods in this cartoon.

we recommend [3]. In this section we only touch on the minimum needed to use the algorithm.

It is important that the argument of  $P_B$  in (32) is not  $z$  but the reflection

$$R_A(z) = 2P_A(z) - z. \quad (33)$$

At a fixed point  $z^* \mapsto z^*$  we have

$$P_B(R_A(z^*)) - P_A(z^*) = 0 \quad (34)$$

which implies

$$P_B(R_A(z^*)) = P_A(z^*) = z_{\text{sol}} \quad (35)$$

is a solution of (31) since it lies in the range of both projections. While this conclusion did not rely on using  $R_A$ , without the reflection the iterate  $z$  does not converge to fixed points.

The RRR iteration rule ensures there is a kind of *local convergence*, independent of whether the problem is feasible. As we will see, this is the only thing the user needs to know. Concerns about global convergence are misplaced since even truly hard problems have formulations (31) with easy projections.

A cartoon of the local behavior is sketched in Figure 6. The trajectory of  $z$  is shown for the limit of small  $\beta$ , where it is a smooth curve, but the convergent behavior holds for any  $\beta \in (0, 2)$ . In both the feasible and infeasible cases there is a pair of *proximal points* on the two sets. These coincide when the problem is feasible. In either case, there is an affine space orthogonal to  $A$  and  $B$  (dashed line) on which  $z$  converges (locally). When the problem is feasible this space is the set of fixed points  $z^*$  of RRR and for any of them

$$\begin{aligned} z_A &= P_A(z^*) \\ z_B &= P_B(R_A(z^*)) \end{aligned}$$

are the coincident solution points. This many-to-one relationship of fixed points to solutions makes RRR different from more familiar fixed-point algorithms, like Newton's root finding method.

The local behavior in the infeasible case is the secret sauce of RRR. The algorithm still generates the points  $z_A$  and  $z_B$ , now from a  $z$  that converges to the orthogonal space while also moving along this space in the direction  $A$  (red) to  $B$  (green). If  $z$  is able to completely converge to the orthogonal space, then the corresponding  $z_A$  and  $z_B$  are a proximal pair. But if the local gap  $\Delta = \|z_B - z_A\|$  is large, the speed of the iterate  $\beta\Delta$  will also be large and in just a few iterations the picture of local constraints will have changed. If by luck that new local picture is feasible, then RRR converges to one of the solution's fixed points. Proximal infeasible points are the bane of the more obvious algorithm that alternates the two projections (and gets stuck on the proximal points).

The likelihood of infeasibility, globally, must be taken seriously when the data that defines  $A$  and  $B$  is noisy (phase retrieval) or the model is insufficiently expressive (machine learning). In this case the best that RRR can offer is a pair of proximal points whose gap  $\Delta$  is exceptionally small, say relative to its starting value. Whether  $z_A$  or  $z_B$  should be used as the best-possible approximate solution depends on the application. The correct choice when training networks is  $z_B$ , since the weights concur over data items in set  $B$ .

### A. Metric auto-tuning

As explained in Section IV, the two projections are just as easy to compute when the metric coefficients  $g_{qi}$  in (27) depend on the BTF  $q$  and the data item  $i$ . The training program in BOOLEARN adjusts these in response to the values of the local gaps

$$\Delta_{qi} = \|z_{qi}^A - z_{qi}^B\|, \quad (36)$$

where  $z_{qi}$  are the trio of  $(w, x, y)$  at BTF  $q$  and data-instantiation  $i$  of the network, and the superscripts denote their projections to  $A$  and  $B$ . The rms value of these local gaps,

$$\Delta_{\text{rms}}^2 = \frac{1}{|\mathcal{H} \cup \mathcal{O}| |\mathcal{D}|} \sum_{q \in \mathcal{H} \cup \mathcal{O}} \sum_{i \in \mathcal{D}} \Delta_{qi}^2 \quad (37)$$

is the same, apart from a contribution from the input nodes, as what we have defined as the gap  $\Delta$ , and whose value is zero at a solution. When  $\Delta_{qi}$  is significantly larger than  $\Delta_{\text{rms}}$ , say when averaged over some number of iterations, it means the variables for BTF  $q$  and data item  $i$  are changing much more from iteration to iteration than the variables for other BTFs and data items. This situation will benefit from an intervention that redistributes the variable changes more evenly, since the solution process in hard problems usually involves some cooperativity. To make variables with a large  $\Delta_{qi}$  change

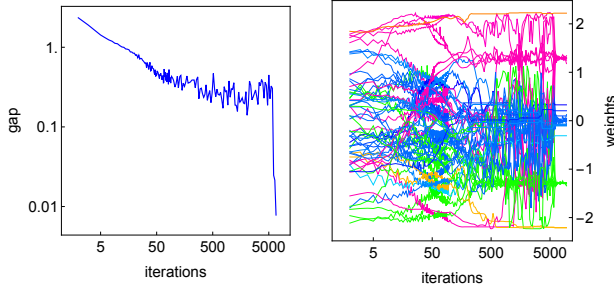


FIG. 7. Evolution of the gap (left) and weights (right) during training of a simple binary multiplier circuit. Though the evolution of the weights is chaotic, the values in the solution (when the gap is very small) are quite simple.

less than their peers, we preferentially increase their metric coefficients by the rule

$$g_{qi} \mapsto g_{qi} + \gamma \left( \Delta_{qi}^2 / \Delta_{\text{rms}}^2 - g_{qi} \right), \quad (38)$$

where  $\gamma > 0$  is the *metric relaxation hyperparameter*. This update is applied in every iteration and preserves the average of the metric coefficients. Initially all metric coefficients are set to the value 1, including those for  $q \in \mathcal{I}$  (which are not included in the update rule). It's important to keep  $\gamma$  small so as not to upset the local convergence of RRR (which holds rigorously only for a constant metric). A good rule of thumb is to set  $\gamma$  near the inverse of the total number of iterations to be performed, so the net change in the metric coefficients can be of order 1.

## VI. APPLICATIONS

### A. Multiplier circuits

Our first application was chosen to make it easy to check the validity of the results. The data is the complete multiplication table, in base 2, for 2-bit integers (from  $0 \times 0 = 0$  to  $3 \times 3 = 9$ ). The network is tasked with learning how the 2+2 bits of the factors are transformed into the 4 bits of the product, by BTFs acting over some number of layers. For maximum interpretability we use  $\sigma = 3$ .

This problem is small enough that through experimentation we are able to trim the depth and width of (fully connected) networks to get a compact representation. The architecture  $4 \rightarrow 4 \rightarrow 4 \rightarrow 4$  appears to be the optimal 3-layer architecture because, empirically, the gap remains positive when the widths of either of the hidden layers is reduced to 3. Figure 7 (left panel) shows the evolution of the gap in a run with hyperparameters  $\beta = 0.2$ ,  $\gamma = 10^{-3}$ . These are good default values, even for much larger problems. The abrupt drop in the gap, enhanced by the logarithmic iterations axis, marks the transition

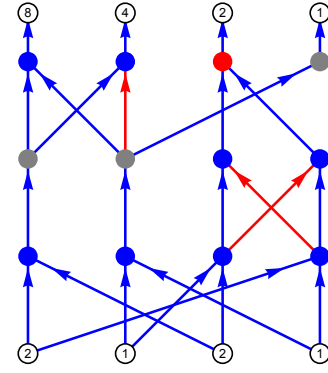


FIG. 8. A three-layer 2-bit multiplier circuit found using BOOLEARN. There are eight AND gates (blue nodes), one OR (red node), and three NOTs (red arrows).

from a period of search to the much faster convergent behavior to the space of fixed points. The corresponding evolution of the 60 weights in the network is shown in the right panel.

As explained at the end of Section III A, the equal-weight-magnitude conclusion of theorem III.4 may not hold when the BTFs do not see all patterns of Boolean values on their relevant inputs in training. Even so, we see that after just a few hundred iterations there is a concentration of weights with magnitude  $\sqrt{5/3} \approx 1.29$ , where 5 is the number of inputs of our BTFs if we include the constant input that implements bias. The evidence suggests that weights are tending toward 3-input gates long before the solution is found. In the short, final convergent phase of the evolution, the conclusion of theorem III.4 is upheld perfectly.

Three of the weights converge to large magnitudes. If a BTF in our network had just a single nonzero weight, its magnitude would be  $\sqrt{5} \approx 2.24$ . The large weights in the solution can have magnitudes smaller than this without changing the 1-relevant-input property of the BTF (the weights to the non-relevant inputs will be small and are also in evidence in the plot).

The circuit found using BOOLEARN is shown in Figure 8. Three of the BTFs (gray nodes) ended up with just one relevant weight and simply copy a BTF output from the layer below. All the others, with three relevant inputs, implement AND (blue nodes) and OR (red node) because one of their relevant inputs is the constant bias node. For the circuit rendering, the signs of the weights incident to hidden nodes were flipped to minimize the number of NOT gates (red arrows). The network's input and output nodes are labeled with the place-values of the bits in the factors and product. A few mathematical facts are brought to light by the circuit. That the 1's bit in the product is the AND of the 1's bits in the factors is the statement that odd numbers always have odd factors. And from the structure of the first two layers we see that multiplication is indeed commutative.

A smaller circuit, with architecture  $4 \rightarrow 5 \rightarrow 4$ , can

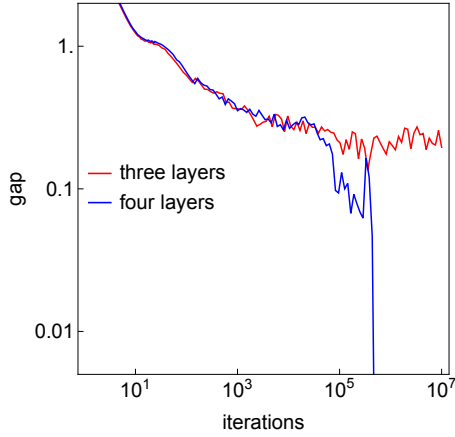


FIG. 9. The gap for learning the 3-bit multiplication table goes to zero for the  $6 \rightarrow 12 \rightarrow 12 \rightarrow 12 \rightarrow 6$  architecture but remains positive when one layer of hidden nodes is removed.

be discovered by using  $\sigma = 5$ . But admitting 5-input majority gates may increase the complexity of the representation. In any case, the exercise is instructive, now because of the three kinds of gates that can arise. In one solution found by BOOLEARN there were three 5-input majority gates (all with bias), five 3-input majority gates (four with and one without bias), and only one BTF with a single relevant input.

Representations of the 3-bit multiplication table (from  $0 \times 0 = 0$  to  $7 \times 7 = 49$ ) should have more layers but also greater width, for storing intermediate results. Figure 9 compares the evolution of the gap for 3- and 4-layer architectures, both having 12 BTFs in each of their hidden layers. With the 4-layer architecture a solution is found in about  $5 \times 10^5$  iterations. In contrast, the gap for the 3-layer architecture appears to fall into a nonzero steady-state, casting doubt on the existence of a solution.

One would not use BOOLEARN to design circuits for  $k$ -bit multipliers! We have no reason to believe the circuits found from data at particular  $k$ 's generalize to arbitrary  $k$ . Even so, this more limited exercise has demonstrated that BOOLEARN generates interpretable representations and that the capacity for depth is not unique to the gradient approach.

## B. Binary encoding-decoding

One of the most convincing demonstrations of learning by error-back-propagation in Rumelhart et al. [4] was the “binary” autoencoder. The authors studied the (fully connected) architecture  $2^k \rightarrow k \rightarrow 2^k$  and gave results for  $k = 3$ . Learning the identity map from data comprising  $2^k$  linearly independent vectors requires nonlinearities in the hidden layer, since otherwise the output can only have the linear span of the narrowest layer ( $k$ ). Rumelhart et al. used 1-hot data vectors and achieved perfect accuracy with the sigmoid activation function. However,

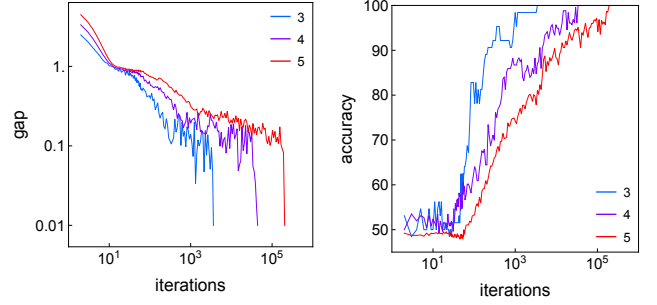


FIG. 10. Evolution of the gaps (left) and training accuracies (right) in the binary autoencoder problem for  $k = 3, 4, 5$ .

their results fell short of true binary encoding-decoding because on some of the data the code values (sigmoid outputs) included the number  $1/2$ .

Unlike the gradient approach, where zero-loss only implies perfect accuracy, when the constraint-based method achieves a gap of zero we know not only that the autoencoder’s accuracy is perfect, but that the code is perfectly binary. Zero-gap implies a feasible point has been found, which includes the property that the concur projections  $y^B$  are equal to the Boolean  $y^A$  projections. Using the same architecture as Rumelhart et al. and 1-hot data vectors (a single  $+1$  replacing one element of a vector of all  $-1$ ), the RRR algorithm easily learns weights for true binary encoding-decoding. We could present those results, not just for  $k = 3$  but also some larger  $k$ 's. Instead, we report on something more interesting.

There is no symmetry between the encoding and decoding stages of the autoencoder, but to see this asymmetry the data vectors have to be more generic than 1-hot vectors. We found that when the data vectors are drawn from the uniform distribution on the hypercube, the simple 2-layer architecture fails, even when the margin constraint is weakened with a large  $\sigma$ . For  $k = 3$  the gap fluctuates indefinitely around the value 0.2. Through experiment we discovered the problem is in the decoder. Augmenting just this stage, so the autoencoder has architecture  $2^k \rightarrow k \rightarrow 2^k \rightarrow 2^k$ , makes the constraint problem feasible.

Figure 10 shows the evolution of the gap and training accuracy with RRR iterations for generic instances of the autoencoder problem and  $k = 3, 4, 5$ . The hyperparameters in these experiments were  $\sigma = 7$ ,  $\beta = 0.2$  and  $\gamma = 10^{-4}$ . The training accuracy is the fraction of output bits that are correct (match the corresponding input bits). This reaches 100% when the gap undergoes a sharp drop. That a gap of 0.01 remains a reliable stopping criterion, even when networks are doubled in size, adds support to our normalization conventions.

The role of the extra processing layer in the decoder is a mystery to us. Figure 11 compares our  $k = 5$  random data vectors with the vectors generated by the first stage ( $k \rightarrow 2^k$ ) of the decoder. Clearly there is something special about the latter, since the 2-layer architecture





FIG. 11. Generic (left) and transformed (right) data vectors (rows) for the  $k = 5$  binary autoencoder problem.

( $2^k \rightarrow k \rightarrow 2^k$ ) would have worked with these as the data. We confirmed this and found that solutions were found on average in only 900 iterations. It would be interesting to know what property allows the transformed data to be decoded in just one layer.

A decoding network can be learned directly in BOOLEARN by using the “any Boolean code” option—constraint (29)—for the inputs. The network now has architecture  $k \rightarrow 2^k \rightarrow 2^k$  and instead of imposing data values at the input nodes in the  $A$  projection, the  $y$  variables are simply rounded to  $-1$  or  $+1$ , whichever is closer. BOOLEARN keeps track of which input codes are getting used, and with what frequency. Since we are training on  $2^k$  random and distinct data vectors in the output, each of the  $2^k$  binary codes must be used exactly once in a solution.

Keeping  $\beta = 0.2$  and  $\gamma = 10^{-4}$ , we were surprised to find that training the decoder, for  $k = 5$ , was much easier with a *smaller* value of  $\sigma$ . Limiting each solution attempt at  $2 \times 10^6$  iterations, the success rate was 100% in 20 attempts for  $\sigma = 5$  (average iteration count  $8.6 \times 10^5$ ), and 0% for  $\sigma = 7$ . The success rate at  $\sigma = 3$  was also 0%, but that is explained if the large BTF margin makes the constraint problem infeasible. On the other hand, any solution found with  $\sigma = 5$  would still be feasible with the smaller margin-bound setting of  $\sigma = 7$ . This finding indicates that tightening the constraints and thereby decreasing the size of the feasible set can make it *easier* for the RRR algorithm to find a point in that set!

### C. Generating MNIST 4’s

Data taken from the wild usually does not come with the guarantee it can be generated by a network of BTFs! To show how BOOLEARN can be used in that setting, we turn to MNIST, a repository of representations of the digits 0–9 used by early humans. We selected just the 4’s, cropped the images to  $16 \times 16$ , and binarized them. The last step is straightforward since the pixel contrast distribution is strongly bimodal. To generate 1024 of the resulting binary strings (each of length  $16^2$ ) with a binary code, as in the last application, would require a 10-bit

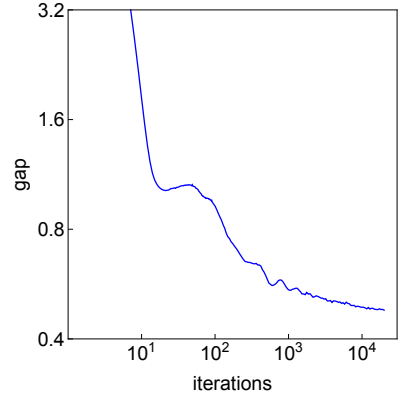


FIG. 12. The gap when training a generator of MNIST 4’s from just 16 1-hot codewords.

code. But there is no reason at all to use the hypercube code for this data, with its large range of distances between codewords. We will instead use the 1-hot code and input-constraint (30). The 1-hot codewords are equidistant, and if we want to try to generate all the data with  $k$  codewords, our network will have  $k$  inputs. Since the data vectors are distinct, our representation of MNIST 4’s will be approximate in some sense because the constraint problem is infeasible. The RRR algorithm finds proximal point-pairs when a problem is infeasible and it will be interesting to see what that means in this application.

Using the same extra-layer decoder design that worked well in the previous application, to build a 16-codeword representation we use a network with architecture  $16 \rightarrow 256 \rightarrow 256$ . Figure 12 shows the evolution of the gap when training on 1024 data (in a single batch) with hyperparameters  $\sigma = 5$ ,  $\beta = 0.2$  and  $\gamma = 10^{-3}$ . The gap has saturated at a value above 0.4. There is clearly no point in performing more iterations.

Recall that in the infeasible case one of the proximal points is on set  $A$ , the other on set  $B$ . In the divide-and-concur setup, the point on  $B$  comes closest to what would be considered an approximate solution because the weights concur across all the data items and the network with those weights is a representation in the usual sense. The companion proximal point, on set  $A$ , has the property that the network’s 1024 instantiations has 1024 outputs, one for each data item. For the point on  $B$  to be close to the point on  $A$ , the 16 data vectors generated by the concurring weights (of point  $B$ , from each of the 16 codewords) try to be close to those 1024 data vectors. We hypothesize that the proximal pair found by RRR corresponds to an optimal partition of the 1024 data into 16 subsets, with the data in each subset close to one of the 16 data generated by the model.

Figure 13 shows the 16 4’s generated by the model. BOOLEARN also outputs how often each codeword was used in the approximate solution. The images in the figure are ranked by decreasing frequency, from 99 times

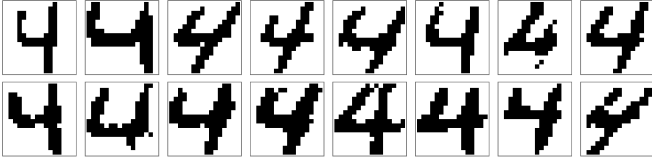


FIG. 13. The 16 archetype 4’s generated by the 1-hot code-word model when trained on 1024 examples. The images are ranked by frequency, from top-left (highest) to bottom right (lowest).

at the top left, to 23 times at the bottom right. The shapes of these “archetype” 4’s are qualitatively reproduced when RRR is rerun with a different random initialization, as are their frequencies. The highly slanted form is always the least frequent.

#### D. Analog data and MNIST

We use the MNIST dataset to show that the discrete BTF outputs are not an impediment to learning analog data. This exercise also provides another example of the power of the margin hyperparameter and letting proximal point-pairs be the endpoint of learning.

To enhance the analog character of the data we use the down-sampled  $8 \times 8$  images available at [11]. Samples of the ten digits are shown in Figure 14. The perceptron model with architecture  $64 \rightarrow 10$  is already interesting in that our method uses a very different approach to learning the weights in the model. In this no-hidden-units network there are just 10 BTFs, one for each digit class. The data constraints (28) impose continuous  $y$  values at the input nodes in the  $A$  constraint and  $\pm 1$  class-encoding values on the  $y$ ’s at the output nodes in the  $B$  constraint. The single correct class node has value  $+1$  and the nine others are set at  $-1$ .

Figure 15 shows the evolution of the gap over  $2 \times 10^4$  iterations when training on 128–2048 items. Since the algorithm is not faced with a combinatorial challenge (multiple near-solutions) we set  $\gamma = 0$ . Also, by using the small time-step  $\beta = 0.05$  there is better convergence to proximal point-pairs—no hurry to get extricated from near-solution traps. The margin hyperparameter, given by  $\mu = \sqrt{65/\sigma}$ , has the expected effect on the gap. A small margin, with  $\sigma = 100$  (left panel), puts a feasible point almost within reach and the gap attains small val-

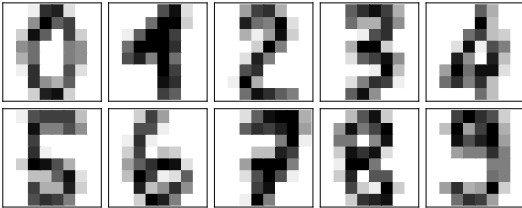


FIG. 14. Down-sampled MNIST digits [11].

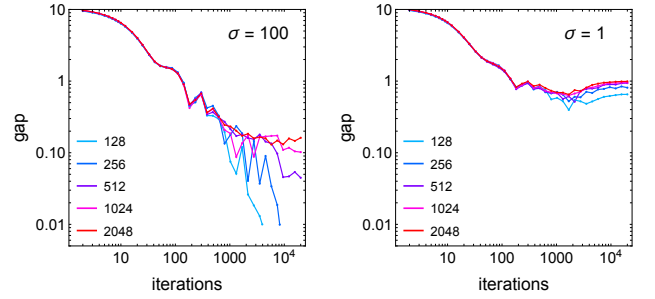


FIG. 15. Effect of the  $\sigma$  hyperparameter on the gap when learning to classify MNIST data for various numbers of training items.

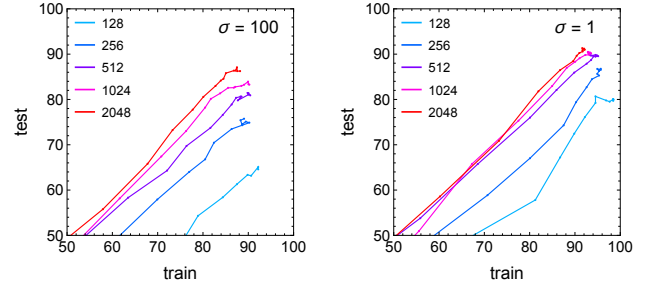


FIG. 16. Detail of the classification-accuracies corresponding to the training runs in Figure 15.

ues. When the margin is large, with  $\sigma = 1$  (right panel), the algorithm is forced to find a proximal point-pair.

The corresponding accuracy plots in Figure 16 make the case that finding a good proximal point-pair is the superior strategy for this kind of data. Recall that the point on the  $B$  constraint defines the solution because the weights concur across all the training items. The class prediction is defined by the output-BTF having the largest value of  $w \cdot y$ , where  $y$  is the vector of grayscale pixel values in a train/test item. Both plots show the learning curves moving toward the main diagonal as the number of training items is increased. But the test accuracy is in general significantly greater for  $\sigma = 1$  than it is for  $\sigma = 100$ . Increasing the BTF margin—by a factor of 10—thereby making the constraint problem infeasible, improved generalization by increasing the separation between the true class and all the wrong classes.

The learned weights, shown in Figure 17, are also interesting. Negative values, shown as red, are just as prominent in the centers of the images as the positive values in blue. Apparently good generalization is just as much about eliminating wrong classes as identifying the correct one. Every detail of Figure 17 is reproduced when the algorithm is rerun from a different random start. The optimal weights for the BTF-perceptron enjoy the same uniqueness, empirically, as the weights obtained in convex optimization.

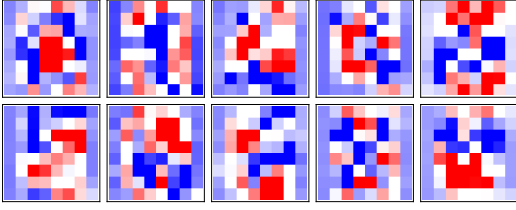


FIG. 17. The learned MNIST weights for the ten digit classes (in the same order as the samples in Figure 14) when training on 2048 items with  $\sigma = 1$ . Blue/red denote positive/negative values.

### E. Logic circuits

We finally come to the kinds of applications where we believe BTF networks can have the greatest impact. The foregoing applications served to contrast the new and existing approaches, while also showing the constraint-based method is not all that different from the perspective of the user. In this section we turn to applications where Boolean representations are natural and have the most to gain from implementations by BTF networks.

Learning Boolean-native data was previewed at the end of Section II. There the BTF networks served as “Boolean inference scratchpads” on which logical rules connecting inputs and outputs were reconstructed. As in any scratchpad, the reconstructions are not unique. For example, logical inferences from the lower layers can simply be copied to higher layers for later use. Given enough layers, arbitrarily complex Boolean functions can be expressed on such scratchpads.

The obvious application that can take advantage of a Boolean representation is language. Words are built from a stem, or root, and various affixed morphemes that modify the meaning and grammatical function. The latter do not form a continuum, but are either absent or present. Clearly a `word2bits` embedding strategy is more natural than continuous embeddings, like `word2vec`. Moreover, if the “atoms” of language are Boolean, then grammar and phrase structure are best represented by Boolean functions. At the highest level is the task of reasoning—the original motivation for the whole Boolean formalism.

The only thing that has kept representations from being Boolean is the reliance of learning algorithms on gradient-based optimizers. For the RRR optimizer, acting on the divide-and-concur variables of a BTF network, discreteness of the representation does not pose a problem. In this section we study just this technical question and leave language modeling for the future.

The circuits shown in Figure 1 implement Boolean functions of the form  $f : \{0, 1\}^m \rightarrow \{0, 1\}^m$  that are compositions of simple “gates” in  $\ell$  layers. BTFs are linearly separable gates, and the  $\sigma$  parameter controls their complexity. We used  $\sigma = 3$  in our experiments, which limits the BTFs to the 1-input COPY gate, 2-input AND/OR, or 3-input MAJ. Our circuits were drawn from uniform distributions so that for large  $m$  two random instances

have the same behavior (and it suffices to study just one). We considered two distributions, both of which used the COPY gate with probability  $1/2$ . The other gates were always AND/OR in one distribution, MAJ in the other. Gate-input nodes from the layer below were drawn from the uniform distribution (of singles, doubles, or triples) with one proviso: no dead-ends in the layer below. This just means the Boolean value at every node gets used by some gate in the layer above. NOT was applied with probability  $1/2$  at every edge (gate-input), including the constant network-input used by the AND/OR gates.

The learning curves plotted in Figure 2 suggest there is a transition to perfect generalization when the number of training items  $N_{\text{data}}$  exceeds some threshold. We can test this hypothesis by comparing the number of circuits with parameters  $(m, \ell, \sigma)$  and the bits of information provided by  $N_{\text{data}}$  items. The number of circuits  $N_{\text{circuit}}$  allowed by the margin constraint with  $\sigma = 3$  is easy to calculate when the gate inputs are sufficiently diverse so that condition (22) holds:

$$N_{\text{circuit}}(m, \ell, 3) = \left( 2^1 \binom{m+1}{1} + 2^3 \binom{m+1}{3} \right)^{\ell m}.$$

The two terms correspond to the  $n = 1$  and  $n = 3$  relevant inputs allowed by  $\sigma = 3$  and the  $m + 1$  node choices include the constant network-input node (used by AND/OR).

The number of distinct Boolean functions  $f$  that can be implemented on a circuit with parameters  $(m, \ell, 3)$  is smaller than  $N_{\text{circuit}}(m, \ell, 3)$  because implementations are far from unique. The greatest source of multiplicity comes from the fact that the nodes within each hidden layer can be freely permuted (without changing edge connectivity) and negation, when applied to all of a hidden BTF’s inputs and outputs, also has no effect on  $f$ . Together, these symmetries give the following lower bound on the multiplicity:

$$N_{\text{mult}}(m, \ell) = (m! 2^m)^{\ell-1}. \quad (39)$$

The number of bits of information required to specify our class of Boolean functions is therefore

$$h(m, \ell, \sigma) = \log_2 N_{\text{circuit}}(m, \ell, \sigma) - \log_2 N_{\text{mult}}(m, \ell).$$

We will treat this as an estimate since the first term is valid only for sufficiently diverse inputs while the multiplicity (39) is only a good lower bound. It evaluates to

$$h(32, 5, 3) = 2466.5 - 598.7 = 1867.8$$

for the parameters in our experiments.

The training algorithm can only acquire these bits of information from the values of the  $m$  outputs bits in each data item. But the information rate is less than  $m$  bits per data item, since these are not uniformly distributed. For example, depending on the circuit, some bits may

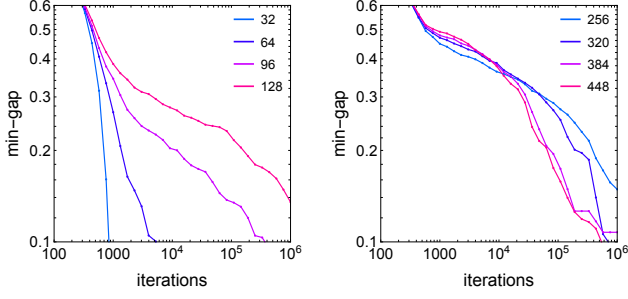


FIG. 18. Minimum value of the gap achieved with increasing iterations for the 5-layer COPY/AND/OR circuit data. On the left is shown the evolution with increasing data in the under-constrained case. As shown on the right, the evolution changes character in the over-constrained case.

be constant or exactly correlated or anticorrelated with other bits. If  $p(x')$  is the probability distribution of output bits  $x' \in \{0,1\}^m$  given by  $x' = f(x)$ , when the inputs  $x$  are uniformly sampled, the information per data item is given by the entropy formula:

$$s(p) = - \sum_{x' \in \{0,1\}^m} p(x') \log_2 p(x') .$$

From the distributions  $p$  produced by feeding  $10^5$  uniform input samples into the two 5-layer, width-32 circuits described above, we obtained

$$\begin{aligned} s(\text{AND/OR}) &= 10.6 \\ s(\text{MAJ}) &= 14.2 . \end{aligned}$$

We then arrive at the following estimates for the number of data items required to specify the two circuits:

$$\begin{aligned} h(32, 5, 3) / s(\text{AND/OR}) &= 176 \\ h(32, 5, 3) / s(\text{MAJ}) &= 132 . \end{aligned}$$

The transition (Figure 2) to apparent perfect generalization falls between 128 and 256 data items for both data sets, consistent with these estimates.

Of course without access to testing data one cannot detect a transition from plots such as Figure 2. Plots of gap vs. iterations can be used instead. Figure 18 shows on the left the evolution of the gap with increasing data (32, 64, 96, 128) when the amount of data is insufficient for perfect generalization. Increasing the data makes the under-constrained problem harder, and more iterations are required to reach a low gap. We have plotted “min-gap”, the lowest currently achieved gap, as this does a better job conveying the evolution. The plots on the right show the evolution when the amount of data is sufficient for perfect generalization (256, 320, 384, 448). Initially (few iterations) the gap is increased with an increase in the data. However, we see that eventually this trend is reversed: increasing the data makes the over-constrained problem easier!

We should not lose track of the fact that good generalization hinges very much on the choice of model. In this application the model is completely specified by just three parameters:  $(m, \ell, \sigma)$ . We believe that decreasing any of these from  $(32, 5, 3)$  makes the constraint problem infeasible with sufficiently many data, and there would be no transition of the kind described above. Moreover, the quality of generalization is independent of the constraint-solver’s hyperparameters, of which there are just two:  $\beta$  and  $\gamma$ . These control the local dynamics of the search (depicted in Figure 6), and not so much the endpoint of the search. In this application we used  $\beta = 0.2$  and  $\gamma = 10^{-3}$  in all the experiments.

## F. Cellular automata

In the Boolean scratchpads we might use to learn some data, there is a tradeoff between width and depth (number of layers). A nice example of this tradeoff is the learning of cellular automata data. Automata data differs qualitatively from random circuit data in that each output bit is a function of just a small number of the input bits. This property may explain our finding that a sparsity-enhancement of the gradient-descent method has a striking positive effect.

Here we consider automata on a 1D periodic world where each cell applies a particular Boolean rule  $R$  to itself and the two cells adjacent:

$$x(p, t+1) = R(x(p-1, t), x(p, t), x(p+1, t)) .$$

The positions  $p$  are periodic with period  $L$  (the size of the world) and  $t$  is the time. Our results are for Wolfram’s “chaotic” rule-30 with formula

$$R_{30}(x, y, z) = x \oplus (y \vee z) ,$$

where  $\oplus$  and  $\vee$  are respectively the exclusive and regular OR operations.

Many of Wolfram’s  $2^8$  rules (on three inputs) can be expressed by a single BTF with negations, but rule-30 is not one of them. Expanding by the two TRUE (or FALSE) cases of the exclusive OR, rule-30 can be expressed in terms of AND/OR as follows:

$$\begin{aligned} R_{30}(x, y, z) &= \text{OR}(\text{AND}(x, -\text{OR}(y, z)), \\ &\quad \text{AND}(-x, \text{OR}(y, z))) \\ &= \text{AND}(-\text{AND}(x, \text{OR}(y, z)), \\ &\quad -\text{AND}(-x, -\text{OR}(y, z))) . \end{aligned} \quad (40)$$

We have switched to  $\pm 1$  for the Boolean values, where negation is multiplication by  $-1$ . This can be expressed on a BTF network with two sets of intermediate values, the first pair holding (a copy of)  $x$  and  $\text{OR}(y, z)$ , and the second pair the two TRUE (or FALSE) cases of the



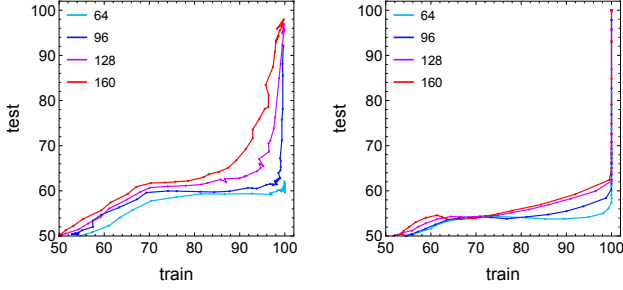


FIG. 19. Evolution of train/test accuracies for learning one step of the rule-30 automaton. The constraint-based results are on the left, gradient-descent with sparsity enhancement is on the right. Both methods used  $10^6$  steps/iterations.

exclusive OR. To implement one time step of rule-30 on a world of size  $L = 16$ , a 3-layer BTF network with architecture

$$16 \rightarrow 32 \rightarrow 32 \rightarrow 16 \quad (42)$$

suffices.

The left panel of Figure 19 shows the evolution of train/test accuracies for 64–160 training data on the architecture (42). It appears that 96 training data are sufficient to learn rule-30 given enough RRR iterations. As in the logic circuit data experiments we used  $\beta = 0.2$  and  $\gamma = 10^{-3}$ . In the right panel are results, for the same number of steps as iterations, of the gradient-descent method. To our surprise, 100% test accuracy was achieved with 96 and 128 training items! The evolution is strange and related to a sparsity enhancement that was also tried on the logic data but without success.

To promote sparsity one can add an  $\ell_1$  penalty on the weights to the model’s loss function. And since each neuron’s input-weights should independently be sparse—like our “support” constraint—we also impose constraint (17). Complete details are given in appendix A. By comparing the evolution of the two kinds of loss (binary-cross-entropy at the outputs, sparsity-penalty throughout the network), it appears the optimizer first manages to

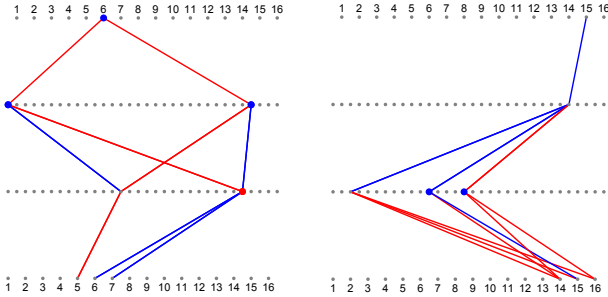


FIG. 20. By thresholding weights by magnitude, to identify the relevant BTF inputs, one can see rule-30 implemented in various ways.

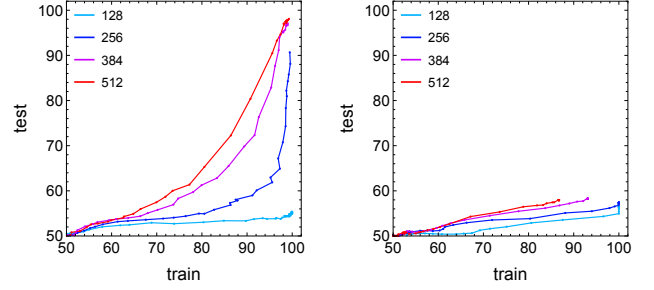


FIG. 21. Same as Figure 19 but for data derived from two steps of rule-30.

perfectly reproduce the correct output bits—without regard to sparsity. In a longer second stage, with training accuracy holding steady at 100%, the sparsity-penalty is reduced as well, until rather abruptly the rule is discovered.

Interpreting the solution-weights to read out the automaton rule is easy in the constraint-based method. This is shown in Figure 20 where, on the left, all the weights above some threshold in magnitude are traced that descend from output bit 6. That includes weights to the network’s constant-input node (not shown). When above the threshold and positive, the node is colored blue and the BTF at that node implements AND. Red nodes have the other sign of weight to the constant-input node and implement OR. The weight-trace shows that output node 6 only depends on input nodes (5, 6, 7), and the Boolean function being implemented exactly matches (40) after the three edges incident to the single COPY node are negated.

Interestingly, from the trace of node 15, shown in the right panel of Figure 20, we see that rule-30 can also be implemented as

$$R_{30}(x, y, z) = \text{MAJ}(-\text{MAJ}(x, y, z), \text{AND}(-x, y), \text{OR}(x, z))$$

after two applications of De Morgan’s rule (to flip the colors of edges). In this implementation the width is greater by one (three nodes instead of two), while the depth is reduced by one (only two layers). The architecture  $16 \rightarrow 48 \rightarrow 16$  could therefore have been used instead. Indeed, learning rule-30 on this architecture is somewhat easier for our method.

Learning rule-30 from input-output patterns separated by two steps of the rule is harder because each bit in the output is now a function of five input bits. From what we just learned about implementing one step of the rule, implementing two steps should be possible with architecture

$$16 \rightarrow 48 \rightarrow 48 \rightarrow 48 \rightarrow 16.$$

Figure 21 shows how the two methods compare on data from two steps of the rule. The iteration/gradient-step limit was doubled relative to the 1-step data on architecture (42), and sparsity enhancement was used as before

in the gradient-descent method. The gradient-descent method is now struggling to even achieve 100% training accuracy.

## VII. NEXT STEPS

Gradient-descent on a loss-objective, and iterated projections to divided BTF constraints, are two strategies for learning with networks. The work for both methods scales in the same way:

$$\frac{\text{computations}}{\text{gradient-step/RRR-iteration}} \propto E \times N_{\text{data}}.$$

Here  $E$  denotes the number of network edges and  $N_{\text{data}}$  the number of data items. The number of gradient steps, or the number RRR iterations, needed for learning depends on the nature of the data, and not much can be said about that. On the other hand, because the constraint approach is able to use BTFs, that strategy has the upper hand when we want the representation to be Boolean.

The advantage of a Boolean representation cannot be overstated. When implemented with BTF networks, and a sufficiently large margin, the BTFs are equivalent to BTFs with very simple weights! As explained in Section III B, even for the generous support parameter  $\sigma = 9$ , the equivalent BTFs all have  $\pm 1$  weights with just one exception. Most of the power consumed by data centers is for inference—e.g. “Are the musicians Justin Bieber and Heinrich Ignaz Biber related?”—and this would be reduced enormously if multiplications are eliminated. Inference in networks with  $\pm 1$  weights uses only integer addition.

Because the constraint-alternative operates at the lowest level of machine learning technology, many of the advances at the higher levels—all trailblazed with gradient-descent—might go through mostly unchanged. Because the representations will become Boolean, there will surely be differences. We hope this article has shown that a core technology based on constraints is at least worth exploring. Below are what we see as the next steps in functionality, implementation, infrastructure, and theory.

### Functionality

BOOLEARN is demonstration software and its programs were deliberately kept simple for that reason. Below are four features that will improve its functionality:

*Weight sharing* This feature is easily implemented in the  $B$  constraint. Weights are always required to concur across data, but additional equalities can be imposed for other reasons, as in convolutional filters.

*Nonuniform margins* There is no reason to impose the same margin constraint on all the BTFs of the network (as in this work). For example, the margins in the

encoder and decoder stages of the autoencoder (Section VIB) should really have been set to different values.

*Batch mode* Whereas the constraint approach to optimization does not rely on “stochasticity” conferred by small data batches, the sheer size of data in many applications makes training—on as many network instantiations as there are data—completely infeasible. A practical solution is to use the largest size batch that can be implemented in hardware, and update the batch in each iteration by replacing the longest residing data item with a new item. With enough iterations all the data ends up getting used, in a cyclic order, where each item’s residence time in the optimizer is equal to the size of the batch.

*Noise accommodation* In phase retrieval, where the data constraints are derived from continuous measurement values (X-ray intensities), solutions are always near-solutions with exceptionally small gaps. The small-gap strategy also works when the data is discrete, and only few bits are corrupted. But there is a better alternative. The trick is to make an allowance for some fraction of the data bits to be flipped from their correct values. In the projection to the data constraint one identifies all the  $y$  variables that are closer to the flipped bit. Of these the allowed fraction are sorted by absolute value and the largest (who change least when the flipped value is indeed the correct one) are projected to the flipped bit. Not only does this make infeasible problems feasible, the corrupted data is flagged in the solution.

### Implementation

The current version of BOOLEARN runs on a single core and all the programs are written in C. Before being considered an alternative to gradient-methods, the following are necessary steps:

*Distributed processing* The independence of the variables associated with each node in both constraint projections, together with the existing low-level C implementation of these projections, makes the framework naturally amenable to parallel and distributed execution. This structural separability has been discussed across multiple formulations and analyses [9, 12–15]. We plan to provide an explicit distributed implementation, and we expect that a variety of parallel and distributed realizations will emerge, tailored to different application domains and hardware environments.

*Python interface* Once distributed execution is in place, the software will be packaged with a Python interface (and potentially interfaces for other languages) to provide a user experience comparable to that of modern gradient-based learning frameworks. In such an interface, the primary optimization diagnostic would be the *gap* rather than a loss value, the RRR time step  $\beta$  plays the role of the learning rate, and the metric relaxation hyperparameter  $\gamma$  is roughly analogous to the moment decay rates of the Adam optimizer. Aside from these dif-

ferences, most user-facing concepts—such as network architectures, batch size, training schedules, and accuracy metrics—would remain closely aligned with standard machine learning practice.

### Infrastructure

We anticipate the development of the following tools to become part of the BTF-learning ecosystem:

*Model interpreters* As explained above, a BTF network trained with a large margin, or equivalently a small  $\sigma$ , is equivalent to a BTF network where all nonzero weights are  $\pm 1$ . Software that inputs a continuous-weight model and outputs an equivalent model with  $\pm 1$  weights is a model interpreter. As explained in connection with the logic circuit reconstructions at the end of Section II, such an interpreter faces a mild complication when the BTF inputs have exact correlations. In that case, interpreters not only have to identify a BTF’s relevant inputs (Section III) but be able to detect correlations in their values.

*Word embedding* As already explained in Section VI E, `word2bits` makes a lot more sense than `word2vec` when the representation is Boolean. Most of the strategies behind `word2vec` carry over to embeddings on the Boolean hypercube.

*Benchmarks for theory* Benchmark data sets are important for keeping machine learning developments in line with the needs of industry. It would be useful to additionally have benchmarks designed to resolve theoretical questions that arise in Boolean representations. How many dimensions are best for embedding the tokens for language-like data? What is the best tradeoff between network width and depth? It’s important that “best” is evaluated not just on the basis of accuracy but also computational efficiency.

### Theory

In their conclusions to “Learning internal representations ...” [4], Rumelhart et al. quote Minsky and Papert on the challenges of establishing “an interesting learning theorem” for multilayered architectures. However, they sum up their own work as follows:

Although our learning results do not *guarantee* that we can find a solution for all solvable problems, our analyses and results have shown that as a practical matter, the error propagation scheme leads to solutions in virtually every case.

This is also where things stand with Boolean threshold function networks and the very different approach we have developed for training them. Despite nonconvexity, the RRR gap is observed to vanish just as consistently,

in diverse applications, as the loss in those gradient experiments 40 years ago. There may never be a rigorous “learning theorem” for either approach, but that should not discourage the development of theoretical analyses that at least make our experience with these methods plausible.

Much work, more than can be reviewed here, has tried to shed light on the gradient-descent process. Many of the same questions can be asked about RRR. Is learning incremental? It certainly seems that way from the quasi-monotone behavior of the gap. The interpretability of the end product of learning with BTF networks makes the study of such questions especially interesting.

### ACKNOWLEDGMENTS

V.E. thanks Chris Myers and Nick Elyu for technical support. M.K.L. acknowledges financial support from the Alexander von Humboldt Foundation and thanks Survir Sra, Adrien Taylor, and Francis Bach for discussions on the future potential of projection methods in deep learning. The authors are grateful to Heinz Bauschke for his influence on their thinking about projection methods.

### DATA AVAILABILITY

Data and code supporting the findings of this study are available at [BOOLEARN](#).

### Appendix A: Gradient-based baselines

Backpropagation-based neural networks have long been used to model logical structure; early work already noted that saturating sigmoids can approximate Boolean behavior [1, 4]. Subsequent work introduced mechanisms for discrete or near-discrete computation within gradient-based training, including straight-through gradient estimators and binary-weight networks [16, 17], as well as differentiable logic-gate architectures that relax discrete choices during training and discretize afterward [18].

This appendix documents the gradient-based baselines used for comparison with BTF networks trained using RRR. The purpose of these baselines is not to develop specialized differentiable circuit learners, but to test whether a standard, fully connected multilayer perceptron trained with widely adopted gradient-based methods can recover the input–output behavior of our Boolean benchmarks. A public leaderboard for the benchmark is available [19].

Our baseline configuration follows common practice in modern deep learning: fully connected ReLU networks trained with AdamW [20], using standard hyperparameters. This setup is consistent with recent empirical studies of Boolean-function learnability that employ comparable architectures and optimization settings [21].

In addition, we tested (i) plain Adam [22] without decoupled weight decay and (ii) a projected-gradient variant enforcing hard row-norm and row-sparsity constraints after each update. In our experiments, neither modification improved test performance relative to the reported baselines, and they are therefore not included in the main comparison.

*Relation to quantization-aware training* Although our method (BTF networks) uses discrete internal values, it is not a variant of quantization-aware training (QAT). In QAT, discretization is introduced during training via surrogate gradients—typically for deployment efficiency—while the underlying optimization remains loss-based and continuous. Representative examples include PACT [23], learned step-size quantization (LSQ) [24], vector-quantized VAEs [25], and post-training quantization methods for transformer models [26–28]. By contrast, our method imposes Boolean node values and hard threshold constraints from the outset, and learning is formulated as a feasibility problem with explicit margin and normalization constraints. Discreteness therefore serves as a structural inductive bias for rule and circuit discovery, not as a compression mechanism.

## 1. Model class and training objective

For each dataset we observe input–output pairs

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^M, \quad x_i \in \{0, 1\}^{d_{\text{in}}}, \quad y_i \in \{0, 1\}^{d_{\text{out}}}.$$

For gradient-based baselines, inputs and outputs are encoded as real values in  $\{0, 1\}$ . We fit a fully connected multilayer perceptron  $f_\theta : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$  with  $L$  affine layers:

$$\begin{aligned} h_0(x) &= x, \\ h_\ell(x) &= \text{ReLU}\left(W^{(\ell)} h_{\ell-1}(x) + b^{(\ell)}\right), \quad \ell = 1, \dots, L-1, \\ f_\theta(x) &= W^{(L)} h_{L-1}(x) + b^{(L)}, \end{aligned}$$

where  $\theta = \{W^{(\ell)}, b^{(\ell)}\}_{\ell=1}^L$  and  $w_j^{(\ell)}$  denotes the  $j$ -th row of  $W^{(\ell)}$ . The final layer is linear and outputs logits; no sigmoid activation is applied inside the network.

The layer widths and depths for each task are specified in the corresponding experimental sections [VIE](#), and [VIF](#). All hidden layers use ReLU activations. No batch normalization, dropout, or other architectural modifications are employed.

Let  $N_{\text{tr}}$  denote the number of training examples. We define the (binary) cross-entropy loss with logits (BCE) as

$$\mathcal{L}_{\text{BCE}}(\theta) = \frac{1}{N_{\text{tr}} d_L} \sum_{i=1}^{N_{\text{tr}}} \sum_{k=1}^{d_L} \left( \log(1 + e^{z_{ik}}) - y_{ik} z_{ik} \right),$$

where  $z_i = f_\theta(x_i)$  are the output logits of the network. This is the average per-output-bit negative log-likelihood over the training set.

## 2. AdamW and Penalty-AdamW

All gradient-based baselines are trained using AdamW [20], i.e. Adam with *decoupled weight decay*. Let  $\theta$  denote the collection of all trainable parameters. At each optimization step  $t$ , we compute the full-batch gradient  $g_t = \nabla_\theta \mathcal{J}(\theta_t)$  on the *entire* training set and apply one AdamW update with learning rate  $\eta$ , moment-decay parameters  $(\beta_1, \beta_2)$ , numerical stabilizer  $\varepsilon$ , and decoupled weight decay coefficient  $\lambda_{\text{decay}}$ , which is applied to all parameters in our implementation.

*AdamW baseline* The AdamW baseline minimizes the data-fit objective

$$\mathcal{J}(\theta) = \mathcal{L}_{\text{BCE}}(\theta),$$

where  $\mathcal{L}_{\text{BCE}}$  is the binary cross-entropy with logits on the training set.

*Penalty-AdamW* To bias optimization toward sparse, circuit-like weight structure, we augment the training objective with a row-norm penalty and an  $\ell_1$  penalty:

$$\begin{aligned} \mathcal{J}(\theta) &= \mathcal{L}_{\text{BCE}}(\theta) + \lambda_{\text{norm}} \sum_{\ell=1}^L \frac{1}{d_\ell} \sum_{j=1}^{d_\ell} \left( \|w_j^{(\ell)}\|_2^2 - d_{\ell-1} \right)^2 \\ &\quad + \lambda_1 \sum_{\ell=1}^L \frac{1}{d_\ell d_{\ell-1}} \sum_{j,k} |W_{jk}^{(\ell)}|. \end{aligned}$$

The row-norm term penalizes deviations from a fixed reference scale, while the  $\ell_1$  term promotes sparsity within rows. Both penalties are applied only to weight matrices (biases are not penalized). No projection step or feasibility enforcement is used; optimization remains entirely gradient-based.

The row-norm penalty fixes a reference scale in an otherwise positively homogeneous ReLU network. Because ReLU satisfies  $\text{ReLU}(cx) = c \text{ReLU}(x)$  for  $c > 0$ , the parameterization is degenerate under layerwise rescalings that leave the network function unchanged. We therefore impose the target  $\|w_j^{(\ell)}\|_2^2 = d_{\ell-1}$  as a fixed, untuned scale convention across layers and datasets. Its role is purely to remove this degeneracy and stabilize optimization. Together with the  $\ell_1$  term, it biases rows toward a small number of large-magnitude entries, promoting sparse, gate-like structure. Although the choice of  $d_{\ell-1}$  parallels the fan-in normalization used in the BTF framework, it carries no Boolean interpretation here and is used only to control weight magnitudes and enable comparison across architectures.

*Hyperparameters and initialization* Unless stated otherwise, all runs use the standard AdamW settings commonly adopted in modern deep learning (and also used in recent work with Boolean data [21]):

$$\eta = 10^{-3} \quad \beta_1 = 0.9 \quad \beta_2 = 0.999$$

$$\varepsilon = 10^{-8} \quad \lambda_{\text{decay}} = 10^{-4}.$$



For **Penalty-AdamW** we additionally fix

$$\lambda_{\text{norm}} = 10^{-2} \quad \lambda_1 = 10^{-4}$$

globally across all datasets and architectures (not tuned per task).

Weights are initialized using the standard Kaiming uniform initialization for ReLU networks [29]. Concretely, for a layer  $\ell$  with fan-in  $d_{\ell-1}$ , weights and biases are sampled independently and uniformly from the interval with bounds  $\pm 1/\sqrt{d_{\ell-1}}$ . This initialization is applied uniformly across all datasets and architectures without modification.

*Batching, step budgets, and logging* All gradient-based baselines use full-batch optimization: each optimization

step consists of one **AdamW** update using the gradient computed on the entire training set. No data subsampling, reshuffling, dropout, or data augmentation is employed.

Random-logic and Rule-30 (1-step) experiments are run for  $10^6$  optimization steps, and Rule-30 (2-step) experiments for  $2 \times 10^6$  steps. We do not use early stopping; all runs are terminated after a fixed number of steps to ensure comparable computational budgets across methods. Loss and accuracies are recorded at checkpoints whose indices are geometrically spaced in iteration number (i.e., the interval between checkpoints increases approximately exponentially). Test-set quantities are computed only for evaluation and never enter the update rule.

- 
- [1] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
  - [2] James R Fienup. Phase retrieval algorithms: a comparison. *Applied Optics*, 21(15):2758–2769, 1982.
  - [3] Veit Elser. *Solving Problems with Projections: From Phase Retrieval to Packing*. Cambridge University Press, 2025.
  - [4] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning internal representations by error propagation. Technical report, University of California, San Diego, La Jolla Institute for Cognitive Science, 1985.
  - [5] Simon Gravel and Veit Elser. Divide and conquer: A general approach to constraint satisfaction. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 78(3): 036706, 2008.
  - [6] Heinz H Bauschke and Jonathan M Borwein. On projection algorithms for solving convex feasibility problems. *SIAM review*, 38(3):367–426, 1996.
  - [7] Ryan O’Donnell and Rocco A Servedio. The Chow parameters problem. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 517–526, 2008.
  - [8] Saburo Muroga, Teiichi Tsuboi, and Charles Richmond Baugh. Enumeration of threshold functions of eight variables. *IEEE Transactions on Computers*, 100(9):818–825, 1970.
  - [9] Veit Elser. Learning without loss. *Fixed Point Theory and Algorithms for Sciences and Engineering*, 2021(1): 12, 2021.
  - [10] Heinz H Bauschke, Manish Krishan Lal, and Xianfu Wang. Projecting onto rectangular hyperbolic paraboloids in Hilbert space. *arXiv preprint arXiv:2206.04878*, 2022.
  - [11] E. Alpaydin and C. Kaynak. Optical recognition of handwritten digits. UCI Machine Learning Repository, 1998.
  - [12] Manish Krishan Lal. *Nonconvex Projections Arising in Bilinear Mathematical Models*. PhD thesis, The University of British Columbia, 2023.
  - [13] Manish Krishan Lal. The flow-limit of reflect-reflect-relax: Existence, stability, and discrete-time behavior. *arXiv preprint arXiv:2512.23843*, 2025.
  - [14] Andreas Bergmeister, Manish Krishan Lal, Stefanie Jegelka, and Suvrit Sra. A projection-based framework for gradient-free and parallel learning. *arXiv preprint arXiv:2506.05878*, 2025.
  - [15] Manish Krishan Lal. Backpropagation from KL projections: Differential and exact I-projection correspondences. *arXiv preprint arXiv:2512.24335*, 2025.
  - [16] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation, 2013. URL <https://arxiv.org/abs/1308.3432>.
  - [17] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in Neural Information Processing Systems*, 28, 2015.
  - [18] Felix Petersen, Christian Borgelt, Thomas Kuehn, and Barbara Hammer. Deep differentiable logic gate networks. In *Advances in Neural Information Processing Systems*, 2022.
  - [19] Manish Krishan Lal. Random majority circuit learning. Kaggle Competition, 2026. URL <https://kaggle.com/competitions/random-majority-circuit-learning>.
  - [20] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
  - [21] Marcio Nicolau, Anderson R. Tavares, Zhiwei Zhang, Pedro H. C. Avelar, João Marcos Flach, Luis DC Lamb, and Moshe Vardi. Understanding boolean function learnability on deep neural networks: PAC learning meets neurosymbolic models. In *19th International Conference on Neurosymbolic Learning and Reasoning*, 2025.
  - [22] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
  - [23] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. PACT: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085*, 2018.
  - [24] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. In *International Conference on Learning Representations*, 2020.
  - [25] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning.

- arXiv preprint arXiv:1711.00937*, 2017.
- [26] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. In *International Conference on Learning Representations (ICLR)*, 2023. arXiv:2210.17323.
  - [27] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning (ICML)*, 2023. arXiv:2211.10438.
  - [28] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023. Extended NeurIPS submission.
  - [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.