

# A parameterized activation function for learning fuzzy logic operations in deep neural networks

Luke B. Godfrey

Department of Computer Science and Computer Engineering  
University of Arkansas  
Fayetteville, AR 72701

Michael S. Gashler

Department of Computer Science and Computer Engineering  
University of Arkansas  
Fayetteville, AR 72701

**Abstract**—We present a deep learning architecture for learning fuzzy logic expressions. Our model uses an innovative, parameterized, differentiable activation function that can learn a number of logical operations by gradient descent. This activation function allows a neural network to determine the relationships between its input variables and provides insight into the logical significance of learned network parameters. We provide a theoretical basis for this parameterization and demonstrate its effectiveness and utility by successfully applying our model to five classification problems from the UCI Machine Learning Repository.

## I. INTRODUCTION

Neural networks are powerful adaptive models with applications in many disciplines. With the backpropagation algorithm, neural networks can be trained in a straightforward manner by gradient-based optimization techniques like stochastic gradient descent and RMSProp [1]. Some of these models, such as convolutional neural networks, learn parameters that can be visualized, interpreted, and understood in some cases [2]. Most neural networks, however, are considered black boxes [3], [4] and it is difficult to determine the semantic meanings behind learned weights.

Fuzzy inference systems, built on fuzzy logic, are also powerful models. Unlike neural networks, fuzzy inference systems are straightforward to interpret and often use linguistic values [5]. In fact, these systems are functionally equivalent to a subset of neural networks [6]. Fuzzy inference systems are less general than neural networks, however, and many neural network techniques are not easily translated into the domain of fuzzy logic.

For these reasons, there is great interest in combining neural networks with fuzzy logic and fuzzy inference systems. Adaptive fuzzy systems have been studied for decades [7], [8], [9], and neural fuzzy modeling continues to be an active topic of research [10], [11], [12]. One purpose of combining these techniques is to produce a model with the flexibility and accuracy of black-box neural networks and the interpretability of fuzzy systems.

Although neural fuzzy systems are well-studied, existing approaches primarily focus on combining specific logical operations in a predefined manner, the most common being an **or** of **ands** [7], [8], [13]. Models that restrict themselves to particular kinds of expressions are limited in the insights they can offer about given datasets. A system that can adaptively

choose from a larger set of logical operations, on the other hand, would be able to provide us with more knowledge about the relationships between its various inputs.

We present a deep learning architecture for learning fuzzy logic expressions by using a novel adaptive transfer function. Our model uses an innovative, parameterized, differentiable activation function that can learn a number of logical operations by gradient descent. Parameters learned by our model can be interpreted as fuzzy rules and combined to form complex logic expressions, allowing a glimpse into the knowledge gleaned during the training process. In Section IV, we report the results of applying our model to five classification problems taken from the UCI Machine Learning Repository [14]. We find that our model is able to learn complex logical expressions and to achieve accuracy comparable to a standard deep neural network with **tanh** activation functions.

## II. FUZZY LOGIC AND NEURAL FUZZY SYSTEMS

Fuzzy logic [15] extends boolean logic into a continuous domain. Typically, **false** is represented as **0**, **true** is represented as **1**, and values in between indicate a corresponding “fuzzy” degree of uncertainty. A typical set of fuzzy operators that generalize the behavior of boolean logic are:

$$\begin{aligned}\mathbf{identity}(x) &= x \\ \mathbf{not}(x) &= 1 - x \\ \mathbf{or}(x, y) &= 1 - (1 - x) \cdot (1 - y) \\ \mathbf{xor}(x, y) &= x + y - 2 \cdot x \cdot y \\ \mathbf{and}(x, y) &= x \cdot y \\ \mathbf{nor}(x, y) &= (1 - x) \cdot (1 - y) \\ \mathbf{nxor}(x, y) &= 1 - (x + y - 2 \cdot x \cdot y) \\ \mathbf{nand}(x, y) &= 1 - x \cdot y\end{aligned}$$

One of the most common applications of fuzzy logic is to control systems [16] using a fuzzy inference system [17]. A typical fuzzy logic controller uses Gaussian membership functions to “fuzzify” inputs (which may be linguistic [5]), a set of (fuzzy) logical IF-THEN rules to apply to the fuzzified inputs, and a function that aggregates and “defuzzifies” the result to a crisp value that determines the system’s action [6]. For example, a fuzzy inference system for an autonomous car might have this inference rule: *IF speed limit IS low OR traffic IS dense THEN speed = slow*. In this example, *low*, *dense*, and *slow* are linguistic values that correspond to functions that map

raw sensor inputs to real numbers in the range [0..1] where 0 indicates definitely not in the given set, 1 indicates definitely in the set, and anything else represents how typical the input is for the given set. 100 km/h might be 0.5 in the moderate speed set and 0.8 in the fast speed set.

The combination of fuzzy logic with neural networks has been termed “fuzzy modeling” [7], “neural fuzzy systems” [9], and “adaptive fuzzy systems” [8]. These systems were extensively studied in the 1990s [18], [19], [20], [21], [22], [23], and one of the primary reasons was that the weights and parameters of a neural fuzzy system could be interpreted more meaningfully than in a traditional neural network [24]. More recently, fuzzy neural networks have been applied to tracking control [11] and other nonlinear dynamics [12].

One kind of neural fuzzy system is any neural network that directly models a fuzzy inference system [25], [26], [19]. These models generally have five layers: 1) an input layer, 2) a membership layer to fuzzify input, 3) a rules layer that computes products of the second layers outputs, 4) a normalization layer to defuzzify signals, and 5) a summation layer to produce the output. Inputs and outputs to this kind of system are crisp, and the fuzzy logic takes place in the hidden layers of the network. In 1993, Jang and Sun showed that these neural networks are functionally equivalent to fuzzy inference systems [6]. The learning that occurs in this kind of neural fuzzy system tunes the member functions, adjusting means and standard deviations, in addition to the combination weights in the output layer. The rules layer in these models is implemented as a product-of-inputs [7], which is a logical **and**. The summation layer performs the logical **or**, and so most of these models result in a logical **or** of **ands** (or a **max** of **mins** [13]).

Another type of neural fuzzy system is any system that uses both fuzzy logic and neural networks as parts of a whole. Chen et. al recently applied this kind of model to solar radiation forecasting, in which the authors used a fuzzy inference system to combine the predictions three separate neural networks like a weighted ensemble [10]. Other models combine fuzzy systems with genetic algorithms [27] or with a Kalman filter [7]. Still others allow inputs, weights, and outputs to be fuzzy [28]. Kwan and Cai proposed the use of “fuzzy neurons” that combine an aggregation function and an activation function with some number of membership functions [13].

### III. APPROACH

We take an approach similar to the common five-layer neural fuzzy system [8], although we use a deeper network. Our model is unique not in its topology in the interpretability of its weights, however, but in the adaptive activation function we use that is able to learn several logical operations. Our activation function is parameterized, continuous, and differentiable, and can therefore be tuned by gradient descent. This has an advantage over existing neural fuzzy systems because it can model more than just an **or** of **ands**, and it has an advantage over other fuzzy neuron approaches because it can be trained by gradient descent instead of set by hand.

The fuzzy logic operators listed in Section II are elegant because they are simple and continuous. However, the symmetry in these equations is difficult to see because our values are not centered about the origin. If we linearly remap these operations by defining **false** to be **-1** instead of **0**, they become:

$$\begin{aligned}
\mathbf{identity}(x) &= x \\
\mathbf{not}(x) &= -x \\
\mathbf{or}(x, y) &= -\frac{(x-1)(y-1)}{2} + 1 \\
\mathbf{xor}(x, y) &= -x \cdot y \\
\mathbf{and}(x, y) &= \frac{(x+1)(y+1)}{2} - 1 \\
\mathbf{nor}(x, y) &= \frac{(x-1)(y-1)}{2} - 1 \\
\mathbf{nxor}(x, y) &= x \cdot y \\
\mathbf{nand}(x, y) &= -\frac{(x+1)(y+1)}{2} + 1
\end{aligned}$$

Then, if we rewrite them in a consistent form, we obtain:

$$\begin{aligned}
\mathbf{identity}(x) &= + \left( \frac{(x+0)(1+0)}{1} - 0 \right) \\
\mathbf{not}(x) &= + \left( \frac{(x+0)(-1+0)}{1} - 0 \right) \\
\mathbf{or}(x, y) &= - \left( \frac{(x-1)(y-1)}{2} - 1 \right) \\
\mathbf{xor}(x, y) &= - \left( \frac{(x+0)(y+0)}{1} - 0 \right) \\
\mathbf{and}(x, y) &= + \left( \frac{(x+1)(y+1)}{2} - 1 \right) \\
\mathbf{nor}(x, y) &= + \left( \frac{(x-1)(y-1)}{2} - 1 \right) \\
\mathbf{nxor}(x, y) &= + \left( \frac{(x+0)(y+0)}{1} - 0 \right) \\
\mathbf{nand}(x, y) &= - \left( \frac{(x+1)(y+1)}{2} - 1 \right)
\end{aligned}$$

In this form, it is much more apparent that there is symmetry that can be leveraged to unify these operations into a single more general operation. There are many possible functions that perfectly express all of these fuzzy logic operations. Three representative solutions are given in Equations 1, 2, and 3.

$$x \textcircled{\alpha} y = \frac{(x + \alpha)(y + \alpha)}{\alpha^2 + 1} - \alpha^2 \quad (1)$$

$$x \textcircled{\alpha} y = \frac{(x + \alpha)(y + \alpha)}{|\alpha| + 1} - |\alpha| \quad (2)$$

$$x \textcircled{\alpha} y = \frac{t}{|t|} \sqrt{|t|} - |\alpha| \quad (3)$$

where  $t = (x + \alpha)(y + \alpha)$

We can plug any of these functions into our table of logical operations to obtain:

$$\begin{aligned}
\mathbf{identity}(x) &= x &= \mathbf{true} \textcircled{0} x \\
\mathbf{not}(x) &= -x &= \mathbf{false} \textcircled{0} x \\
\mathbf{or}(x, y) &= \mathbf{not}(x \textcircled{1} y) &= \mathbf{false} \textcircled{0}(x \textcircled{1} y) \\
\mathbf{xor}(x, y) &= \mathbf{not}(x \textcircled{0} y) &= \mathbf{false} \textcircled{0}(x \textcircled{0} y) \\
\mathbf{and}(x, y) &= x \textcircled{1} y \\
\mathbf{nor}(x, y) &= x \textcircled{1} y \\
\mathbf{nxor}(x, y) &= x \textcircled{0} y \\
\mathbf{nand}(x, y) &= \mathbf{not}(x \textcircled{1} y) &= \mathbf{false} \textcircled{0}(x \textcircled{1} y)
\end{aligned}$$

Figure III shows a plot of Equation 2 with 5 values for  $\alpha$ .

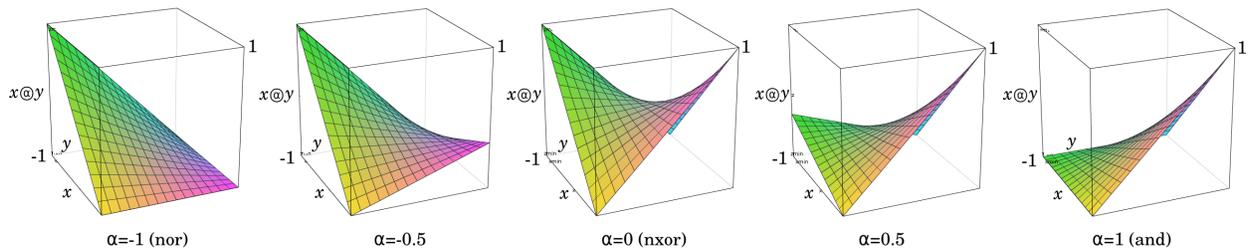


Fig. 1. Equation 2 continuously interpolates among three fuzzy logic operations: **nor**, **nxor**, and **and**. By allowing biases (**true** and **false**) and weights (in particular, a weight of -1), this equation can also compute **identity**, **not**, **or**, **xor**, and **nand**.

### A. Learning Simple Logic Operations

Since Equations 1, 2, and 3 are continuous and differentiable, gradient-based optimization techniques could potentially be used with them to find the values for  $\alpha$  that approximate the logic represented in a set of training examples.

An important consideration in using gradient descent with fuzzy logic that does not typically occur in more traditional applications for gradient descent is that each training example only provides information about a subset of the parameter space. For example, suppose  $\alpha$  is initialized to a random value between  $-1$  and  $1$ , and suppose the training pattern  $1 \otimes 1 = 1$  is presented for optimizing the value of  $\alpha$  by gradient descent. This training pattern suggests that  $\alpha$  should not be less than  $0$ , because  $(1 \text{ nor } 1) \neq 1$ . However, this training pattern does not suggest anything about what specific value  $\alpha$  should take  $\geq 0$ , because  $(1 \text{ nxor } 1) = 1$  and  $(1 \text{ and } 1) = 1$  are both equally true.

Figure III-A shows a comparison of Equations 1, 2, and 3 for the case of computing  $1 \otimes 1$ . If  $\alpha$  has a value less than  $0$ , then gradient-based optimization methods will adjust  $\alpha$  by moving it closer to  $0$  no matter which of these three equations is used. As long as the curve in this region is monotonic, the precise shape is not important because these values are not defined in boolean logic. However, if Equation 1 is used, and  $\alpha$  has a value greater than  $0$ , then gradient-based optimization methods will move  $\alpha$  closer to  $0$  or  $1$ , whichever is closer to the current value of  $\alpha$ . This is incorrect behavior because this training pattern does not provide any information about whether the logical operation should be more like **nxor** or more like **and**. Since both of these operations are consistent with the training pattern, it would be arbitrary to bias the model in favor of one over the other. Arbitrary parameter adjustments are likely to fight against subsequent training pattern presentations that may convey valid information for that region of the parameter space, resulting in the model getting stuck in a local optimum. Equations 2 and 3 correctly adjust  $\alpha$  in all regions of the parameter space with this training pattern.

As another representative case, consider the training pattern  $1 \otimes -1 = 1$ . Perhaps counterintuitively, this pattern provides no information about any region of the parameter space, because  $1 \text{ nor } -1$ ,  $1 \text{ nxor } -1$ , and  $1 \text{ and } -1$  all evaluate to  $-1$ . Correct behavior, therefore, should not adjust the value of  $\alpha$ ,

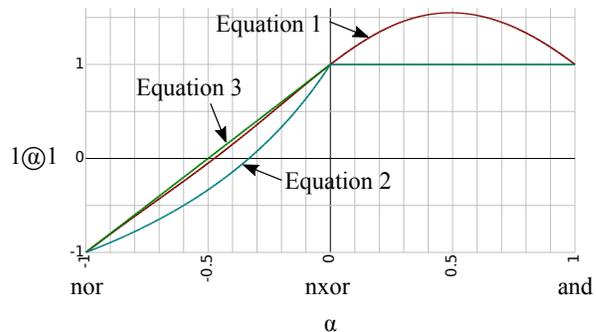


Fig. 2. A comparison of Equations 1, 2, and 3 when computing **true@true**. Equation 1 is smoother, but Equation 2 is better suited for use with gradient-based optimization.

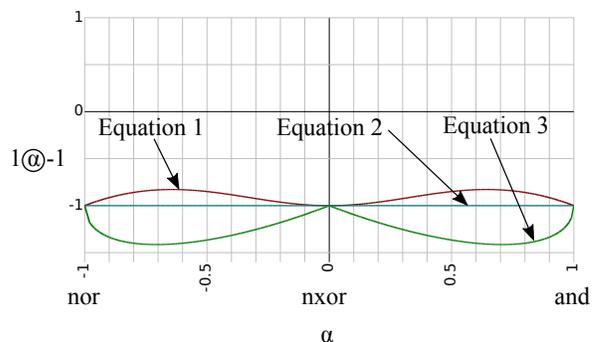


Fig. 3. A comparison of Equations 1, 2, and 3 when computing **true@false**. Only Equation 2 is well-suited for use with gradient-based optimization.

regardless of its current value. Figure III-A shows that only Equation 2 exhibits the correct behavior for this case. (This does not imply that such patterns should be discarded because in a network containing many fuzzy logic units, different input values would reach each of the units depending on the current  $\alpha$  values.)

The remaining two cases are both mirror images of these cases, so the same analysis applies. It follows that Equation 2 can be expected to yield correct behavior with gradient-based optimization in all cases when training with pure boolean

logic. Consistent with this intuition, we found experimentally that Equation 2 was always able to learning simple boolean logical expressions, while the other equations sometimes became stuck in local optima. Therefore, we used Equation 2 to implement  $x@y$  with the remainder of our experiments.

At this point, we must note a potential problem with our approach: Equation 2 is not a t-norm. A t-norm must be commutative, monotonic, and associative, and the value 1 must be the identity element [29]. Although our equation is commutative and monotonic, it is not associative and the value 1 is not neutral for all  $\alpha$ . Our work is, therefore, not t-norm fuzzy logic but belongs instead to a broader class of fuzzy logic. For the purpose of this paper, we use a relaxed definition of fuzzy logic as being any logic with continuous values between **true** and **false**.

### B. Learning Complex Logic Expressions

We refer to a layer of network units that implement Equation 2 as a *Fuzzy layer*. Thus, Equation 2 is used as a sort of adaptive transfer function. The only parameters to train in a fuzzy layer are one  $\alpha$  value per each unit. Equations 4, 5, and 6 give the partial derivatives of Equation 2, which are necessary to train such layers with gradient-based optimization methods.

$$\frac{\partial @}{\partial x} = \frac{y + \alpha}{|\alpha| + 1} \quad (4)$$

$$\frac{\partial @}{\partial y} = \frac{x + \alpha}{|\alpha| + 1} \quad (5)$$

$$\frac{\partial @}{\partial \alpha} = \frac{|a|(x + y) - a(xy + 1)}{|a|(|a| + 1)^2} \quad (6)$$

Most gradient-based optimization methods can be implemented in three steps: (1) A forward propagation step that computes predictions with current values, (2) a backpropagation step that computes “blame” terms for each layer in the model, and (3) an update step that refines the parameters of the model. Equation 2 is used in the forward propagation step. Equations 4 and 5 are used in the backpropagation step to assign blame to preceding layers. Equation 6 is used in the update step to refine the  $\alpha$  values.

Because Equation 6 is not continuous at  $\alpha = 0$ , special care must be taken in the implementation of Fuzzy layers to ensure that they do not become stuck at this point. We addressed this problem in our implementation by adding the statement “**if**  $\alpha < \epsilon$  **then**  $\alpha \leftarrow -\alpha$ ” to our update step. This small addition enables  $\alpha$  to cross over the value 0 in cases where it would otherwise become stuck. As long as  $\epsilon$  is a small value, this will have negligible impact on training precision. We used the value  $\epsilon = 0.001$ .

Another challenge that arises in learning fuzzy logic is that Equation 2 only accepts two input values,  $x$  and  $y$ . One possible solution is to try to generalize the equations in a manner that can support vectors of arbitrary dimensionality. Equation 3 can be generalized in this manner, as given in Equation 7.

$$f(\vec{x}, \alpha) = \frac{t}{|t|} |t|^{1/n} - |\alpha| \quad (7)$$

$$\text{where } t = \prod_i^n (x_i + \alpha)$$

(In Equation 7,  $f$  is the fuzzy operator, and  $n$  refers to the number of elements in  $\vec{x}$ .) In higher dimensions, **or** becomes **any**, **and** becomes **all**, and **xor** becomes **parity**. Unfortunately, Equation 3 resists gradient-based optimization, so its generalized version is unlikely to do any better, and Equation 2 cannot be generalized in this manner. Further, in applications with many variables, it is often unlikely that all variables will simultaneously take the same value, which renders the **any** and **all** operations to have very limited utility. Therefore, a good solution should provide a mechanism to select which values feed into each logical operation. This is also consistent with most real-world uses of boolean logical expressions, and logical expressions involving only two variables at a time are more likely to be easily comprehensible to humans than logic involving many values. Our solution to this challenge introduces two additional layer types, which we call *AllPairings*, and *FeatureSelector*.

An AllPairings layer accepts  $n$  inputs, and outputs all  $n(n-1)/2$  possible unordered pairings of its input values. In order to facilitate the **identity** and **not** operations, we additionally pair each input value with the bias values **true** and **false**, which increases the total number of unordered output pairs to  $n(n-1)/2 + 2n$ . This layer type contains no parameters to train, so it is straightforward to use in a deep network. During the backpropagation step, the blame term assigned to each input unit is simply the sum of the blame terms for all affect output units.

A FeatureSelector layer is identical to a traditional fully-connected linear layer, except with four minor modifications: (1) No bias weights are used, (2) The weights are initialized with uniform values instead of random values, (3) The weights that feed into each unit are constrained to have values between -1 and 1, and (4)  $L^1$  regularization is applied to the weights in this layer to gently promote sparse connections in this layer. As it is closely related to a fully-connected layer, it produces a weighted sum of its inputs. This allows the network to learn an interpolation between logical expressions learned in the Fuzzy layers. The outputs of a FeatureSelector layer can be remapped between -1 and 1 (i.e. through a membership function) before they are used as input to other layers.

Our topology, given the definitions of these two layer types, is as follows. Given  $n$  continuous input values, we first normalize them between -1 and 1; this can be thought of as a single linear membership function (where 1 is “high” and -1 is “low”). Next, we feed the normalized values into an AllPairings layer. We then feed all combinations of value pairings into a Fuzzy layer, which learns an optimal logical operation for each pair of values. The output of the Fuzzy layer is fed into a FeatureSelector layer to manage dimensionality and to produce the desired number of output values. If a deeper topology is desired, we can feed the output of the FeatureSelector layer into another normalizing layer (i.e. a

TABLE I  
AVERAGE ERRORS OF A DEEP NEURAL NETWORK (DNN), AN ENSEMBLE OF ANFIS FUZZY CLASSIFIERS (EFC), OUR MODEL, AND THE EXPRESSIONS OBTAINED FROM “SNAPPING” THE WEIGHTS IN OUR NETWORK (SNAPPED) ON THE VALIDATION DATA FOR FIVE CLASSIFICATION PROBLEMS. BEST RESULTS ARE **BOLD**ED.

Dataset	DNN	EFC	Our Model	Snapped
Breast Cancer	3.26%	6.42%	<b>2.77%</b>	2.84%
Diabetes	29.68%	24.51%	<b>22.79%</b>	35.06%
Vehicle	<b>18.01%</b>	50.58%	28.71%	67.84%
Waveform	<b>14.95%</b>	-	15.27%	68.43%
Yeast	<b>46.12%</b>	67.37%	49.77%	82.94%

single membership function for each output) and repeat the sequence of AllPairings, Fuzzy, and FeatureSelector layers to an arbitrary depth. If no further depth is needed, we can feed the output of the final FeatureSelector layer into any kind of output layer; we use a **max** layer for classification.

In our validation, we fix the depth of logic layer to two, resulting in the following final 10-layer deep topology: (1) the input layer (identity), (2) a normalization layer (a single linear membership function), (3) an AllPairings layer, (4) a Fuzzy layer, (5) a FeatureSelector layer, (6) a **tanh** layer (a single nonlinear membership function with a new input space), (7) another AllPairings layer, (8) another Fuzzy layer, (9) another FeatureSelector layer, and (10) a **max** layer for classification.

#### IV. VALIDATION

Our goal in this work is not to surpass the accuracy of existing results or to claim any novelty in learning fuzzy logic rules, but to demonstrate an alternative approach to fuzzy learning using a novel adaptive transfer function. We apply our model to five classification problems taken from the UCI Machine Learning Repository [14], comparing it with a regular deep neural network (DNN) with **tanh** activation functions as well as with ensemble of fuzzy classifiers (EFC) proposed by Canul-Reich et. al in 2007 [30]. Our results demonstrate that our model meets our goal, adaptively learning complex logic expressions by gradient descent while yielding accuracies comparable to existing methods and learning fuzzy logic rules.

The five classification problems we used were breast cancer, diabetes, vehicle, waveform, and yeast. Errors yielded by four models are compared in Table I. Across all five problems, we fixed our models parameters to demonstrate robustness. In particular, we set the learning rate to 0.01 and the regularization term to 0.0001. We use a DNN topology with a depth and number of weights similar to our model, but with each layer being fully connected and all activation functions set to **tanh**.

We also include the results reported by Canul-Reich et. al [30] rather than re-evaluating their method ourselves (which is why there is no result for their model on the waveform problem). Their model is an ANFIS-based ensemble of fuzzy classifiers (labeled EFC in Table I). One of the important observations made by Canul-Reich et. al was that their model performed poorly on datasets with more than six input features (vehicle and yeast).

All of the datasets we use have continuous inputs and discrete class outputs. Breast cancer has 9 inputs and 2-class outputs. Our model outperforms both DNN and EFC in terms of accuracy. Diabetes (Pima) has 8 inputs and 2-class outputs. Again, our model achieves higher accuracy than the compared models.

Vehicle has 18 inputs and 4-class outputs. This is an interesting problem because it has a high number of inputs and more than a binary output. The DNN has significantly lower error than our model, likely due to the difficulties fuzzy systems typically have with larger input spaces. However, our model makes more accurate predictions than EFC, demonstrating that although our model is susceptible to problems with large input spaces, it is more robust than some other fuzzy-based approaches.

Waveform has 40 inputs and 3-class outputs. This is a very large input space, so it is an interesting test for a fuzzy-based system like ours. Our model performed well on this problem, yielding test errors only marginally higher than the DNN. (Canul-Reich et. al did not report results on this problem for EFC).

Yeast has 8 inputs and 10-class outputs. Out of the five datasets we used, this had the highest number of output classes, so the results of this test demonstrate how our model handles many classes. Our model performed nearly as well as the DNN and significantly better than the EFC, demonstrating that it is able to model data with many partitions.

Our model’s weights can be interpreted as complex logical expressions that describe the relationships learned between various inputs. We form these expressions by “snapping” the  $\alpha$  parameters to the nearest whole value (i.e.  $\alpha = -0.2$  snaps to 0 to become **nxor**) and “snapping” the weights on the FeatureSelector layer to the nearest whole value (effectively negating expressions with weights near -1 and dropping expressions with weights near 0). We use  $n$  to denote that input  $n$  is “high”,  $\neg$  for **not**,  $\&$  for **and**,  $|$  for **or**, and  $\oplus$  for **xor**.  $c_i$  is the output for class  $i$ , where the  $i$  with maximum  $c_i$  is the predicted class. Addition has no equivalent logical operation, so we leave it as a sum indicating the interpolation between operands. (In the network, the sum is fed through a **tanh** function to map the result back into our logical space; this step is omitted in the snapped expressions). The resulting expressions for three of the problems (breast cancer, diabetes, and vehicle) are shown in Table II; expressions for waveform and yeast were omitted from this paper for brevity. Accuracies of the formed expressions applied to the validation data are reported in the “Snapped” column of Table I.

#### V. CONCLUSION

We presented a deep learning architecture for learning fuzzy logic expressions. The primary component of this architecture is an innovative, parameterized, differentiable activation function that can learn a number of logical operations by gradient descent. This activation function is unique to our approach and unifies fuzzy logic with deep learning in a new way that leverages the advantages of both domains. We provided

TABLE II

LOGICAL EXPRESSIONS LEARNED BY OUR MODEL FOR THREE OF THE DATASETS. THESE EXPRESSIONS WERE FORMED BY “SNAPPING” ALL PARAMETERS OF THE NETWORK TO THE NEAREST WHOLE VALUE. EXPRESSIONS FOR THE OTHER TWO DATASETS WERE OMITTED FOR THE SAKE OF BREVITY.

Dataset	Expression
Breast Cancer	$c_0 = (0   3) + (1   5)$ $c_1 = ((1   3) + (3   7))   ((2   8) + (3   4) + (5   6))$
Diabetes	$c_0 = (\neg(2   4) + 7 + \neg(2   6))   (\neg(2 \& 6) + \neg 0)$ $c_1 = (\neg(2 + \neg(1 \& 5) + 1) \& (2   6))$
Vehicle	$c_0 = \neg((0 \& 4)   ((6 \& 10) + (4 \& 16) + (7 \& 13)))$ $\quad + (\neg(2   16) \& ((4 \& 9) + (5 \& 6)))$ $c_1 = (4 \& 9) + (5 \& 6) + \neg(((0 \& 4) \oplus (8 \& 16)))$ $c_2 = \neg(((3 \& 5) + (0 \& 4) + (3 \& 13)) \& (0 \& 4))$ $c_3 = \neg((5 \& 15) \& (0 \& 4))$

both a theoretical basis for our activation function and testing results on five classification problems from the UCI Machine Learning Repository. Not only was our model a reasonably strong classifier for these problems, but also the parameters of our model were interpretable as logical expressions that summarized what it had learned. Our primary contribution is neither a better classifier nor the first model whose weights can be interpreted as fuzzy logic rules, however, but the adaptive transfer function that learns logical operations by gradient descent. We believe that this is a promising new approach to combining fuzzy logic with deep learning that potentially opens the way for future work in both domains.

## REFERENCES

- [1] T. Tieleman and G. Hinton, “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude,” *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, 2012.
- [2] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks,” in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [3] D. Kanungo, M. Arora, S. Sarkar, and R. Gupta, “A comparative study of conventional, ann black box, fuzzy and combined neural and fuzzy weighting procedures for landslide susceptibility zonation in darjeeling himalayas,” *Engineering Geology*, vol. 85, no. 3, pp. 347–366, 2006.
- [4] J. Sjöberg, Q. Zhang, L. Ljung, A. Benveniste, B. Delyon, P.-Y. Glorenec, H. Hjalmarsson, and A. Juditsky, “Nonlinear black-box modeling in system identification: a unified overview,” *Automatica*, vol. 31, no. 12, pp. 1691–1724, 1995.
- [5] L. A. Zadeh, “Fuzzy logic= computing with words,” *IEEE transactions on fuzzy systems*, vol. 4, no. 2, pp. 103–111, 1996.
- [6] J.-S. Jang and C.-T. Sun, “Functional equivalence between radial basis function networks and fuzzy inference systems,” *IEEE transactions on Neural Networks*, vol. 4, no. 1, pp. 156–159, 1993.
- [7] J.-S. R. Jang et al., “Fuzzy modeling using generalized neural networks and kalman filter algorithm,” in *AAAI*, vol. 91, 1991, pp. 762–767.
- [8] J.-S. Jang, “Anfis: adaptive-network-based fuzzy inference system,” *IEEE transactions on systems, man, and cybernetics*, vol. 23, no. 3, pp. 665–685, 1993.
- [9] C.-T. Lin and C. G. Lee, “Neural fuzzy systems,” *PTR Prentice Hall*, 1996.
- [10] S. Chen, H. Gooi, and M. Wang, “Solar radiation forecast based on fuzzy logic and neural networks,” *Renewable Energy*, vol. 60, pp. 195–201, 2013.
- [11] C. P. Chen, Y.-J. Liu, and G.-X. Wen, “Fuzzy neural network-based adaptive control for a class of uncertain nonlinear stochastic systems,” *IEEE Transactions on Cybernetics*, vol. 44, no. 5, pp. 583–593, 2014.
- [12] E. Kayacan, E. Kayacan, and M. A. Khanesar, “Identification of nonlinear dynamic systems using type-2 fuzzy neural networks a novel learning algorithm and a comparative study,” *IEEE Transactions on Industrial Electronics*, vol. 62, no. 3, pp. 1716–1724, 2015.
- [13] H. K. Kwan and Y. Cai, “A fuzzy neural network and its application to pattern recognition,” *IEEE transactions on Fuzzy Systems*, vol. 2, no. 3, pp. 185–193, 1994.
- [14] A. Asuncion and D. J. Newman, “UCI machine learning repository,” 2007. [Online]. Available: <http://archive.ics.uci.edu/ml/>
- [15] G. Klir and B. Yuan, *Fuzzy sets and fuzzy logic*. Prentice hall New Jersey, 1995, vol. 4.
- [16] C.-C. Lee, “Fuzzy logic in control systems: fuzzy logic controller. i,” *IEEE Transactions on systems, man, and cybernetics*, vol. 20, no. 2, pp. 404–418, 1990.
- [17] L. A. Zadeh, “Fuzzy logic,” *Computer*, vol. 21, no. 4, pp. 83–93, 1988.
- [18] B. Kosko, “Neural networks and fuzzy systems: a dynamical systems approach to machine intelligence/book and disk,” *Vol. 1 Prentice hall*, 1992.
- [19] Y.-C. Chen and C.-C. Teng, “A model reference control structure using a fuzzy neural network,” *Fuzzy sets and Systems*, vol. 73, no. 3, pp. 291–312, 1995.
- [20] Y. Hayashi, J. J. Buckley, and E. Czogala, “Fuzzy neural network with fuzzy signals and weights,” *International Journal of Intelligent Systems*, vol. 8, no. 4, pp. 527–537, 1993.
- [21] G. A. Carpenter, S. Grossberg, N. Markuzon, J. H. Reynolds, and D. B. Rosen, “Fuzzy artmap: A neural network architecture for incremental supervised learning of analog multidimensional maps,” *IEEE Transactions on neural networks*, vol. 3, no. 5, pp. 698–713, 1992.
- [22] H. Zhenya, W. Chengjian, Y. Luxi, G. Xiqi, Y. Susu, R. C. Eberhart, and Y. Shi, “Extracting rules from fuzzy neural network by particle swarm optimisation,” in *Evolutionary Computation Proceedings, 1998. IEEE World Congress on Computational Intelligence., The 1998 IEEE International Conference on*. IEEE, 1998, pp. 74–77.
- [23] J.-S. R. Jang, C.-T. Sun, and E. Mizutani, “Neuro-fuzzy and soft computing, a computational approach to learning and machine intelligence,” 1997.
- [24] E. Cox, “Adaptive fuzzy systems,” *IEEE Spectrum*, vol. 30, no. 2, pp. 27–31, 1993.
- [25] C.-T. Lin and C. S. G. Lee, “Neural-network-based fuzzy logic control and decision system,” *IEEE Transactions on computers*, vol. 40, no. 12, pp. 1320–1336, 1991.
- [26] C.-T. Lin and C. G. Lee, “Reinforcement structure/parameter learning for neural-network-based fuzzy logic control systems,” *IEEE Transactions on Fuzzy Systems*, vol. 2, no. 1, pp. 46–63, 1994.
- [27] F. Valdez, P. Melin, and O. Castillo, “An improved evolutionary method with fuzzy logic for combining particle swarm optimization and genetic algorithms,” *Applied Soft Computing*, vol. 11, no. 2, pp. 2625–2632, 2011.
- [28] R. J. Kuo, C. Chen, and Y. Hwang, “An intelligent stock trading decision support system through integration of genetic algorithm based fuzzy neural network and artificial neural network,” *Fuzzy sets and systems*, vol. 118, no. 1, pp. 21–45, 2001.
- [29] M. M. Gupta and J. Qi, “Theory of t-norms and fuzzy inference methods,” *Fuzzy sets and systems*, vol. 40, no. 3, pp. 431–450, 1991.
- [30] J. Canul-Reich, L. Shoemaker, and L. O. Hall, “Ensembles of fuzzy classifiers,” in *Fuzzy Systems Conference, 2007. FUZZ-IEEE 2007. IEEE International*. IEEE, 2007, pp. 1–6.