

# A CONCRETE EXAMPLE: BOOLEAN VARIANT

WEIFAN CHEN

ABSTRACT. A very plain example of how to conduct forward path (FP) and backward path (BP) using boolean variant.

## 1. NOTATION

All indices are implicitly counted from 1, because sometimes zero is reserved for other meanings.

$\mathbb{B} = \{T, F\}$ : the boolean set.  
 $\mathbb{N}$ : natural number set.

$x_{k,i}^l \in \mathbb{B}$ : this variable belongs to the  $k^{\text{th}}$  sample in the batch. It is the  $i^{\text{th}}$  input of the  $l^{\text{th}}$  layer. It is also the output of the  $(l-1)^{\text{th}}$  layer. When  $l = 1$ , this indicates the sample input.

$w_{i,j}^l \in \mathbb{B}$ : the weight associated to  $l^{\text{th}}$  layer input. It is the weight that operates between the  $i^{\text{th}}$  output of the  $l^{\text{th}}$  layer and the  $j^{\text{th}}$  input of the downstream unit. When  $i = 0$ , this indicates the bias to the  $j^{\text{th}}$  input of the downstream unit.

$L \in \mathcal{F}(\mathbb{B}, \mathbb{B})$ : a boolean-input boolean-valued function. We use XNOR throughout this example. Thus  $L = \text{XNOR}$  from now on.

$s_{k,j}^l \in \mathbb{N}$ : the pre-activated sum belongs to the  $l^{\text{th}}$  layer. It is the  $j^{\text{th}}$  pre-activated value. It is calculated as:

$$s_{k,j}^l = w_{0,j}^l + \sum_{i=1}^M L(w_{i,j}^l, x_{k,i}^l) \quad (1)$$

Note: the bias here doesn't make sense to me as it's really supposed to be a constant, it's just a boolean value? I remember looking at the paper, and I also had the same issue there.

Then the next layer input can be calculated as:

$$x_{k,j}^{l+1} = \sigma_\tau(s_{k,j}^l) \quad (2)$$

in which  $\sigma_\tau \in \mathcal{F}(\mathbb{N}, \mathbb{B})$  is a natural-input boolean-value activation function defined as:

$$\sigma_\tau(x) = \begin{cases} T & \text{if } x \geq \tau \\ F & \text{if } x < \tau \end{cases} \quad (3)$$

$\tau$  is taken to be constant in this example. It can be half of the number of next layer inputs.

$\text{loss}_k \in \mathcal{F}(\mathbb{B}, \mathbb{B})$ : the boolean-input boolean-value function that maps to the loss for

the  $k^{\text{th}}$  sample in the batch. Throughout this writing, it is XOR for indicating the mismatch between the prediction ( $x_k$ ), and the ground truth label ( $y_k$ ). Formally:

$$\text{loss}_k(x_k) = \mathbf{xor}(y_k, x_k) \quad (4)$$

$\mathcal{L}$ : the accumulated loss function defined as:

$$\mathcal{L} = \sum_{i=1}^K \text{loss}_k(x_k) \quad (5)$$

in which  $K$  is the number of samples in the batch.

## 2. IMPORTANT CONCEPTS

### 2.1. Variant and derivative.

The notation  $f'(x)$  is commonly used for real-number input function. To operate on a boolean DNN, the notation has to adapt for natural-number input functions and boolean-input functions. To do so, first define *variant*  $\delta$  and its related operations.

Definition: variant

For **any** value  $a$  and  $b$ , the variant from  $a$  to  $b$  is:

$$\delta(a \rightarrow b) = \begin{cases} T & \text{if } b > a \\ 0 & \text{if } a = b \\ F & \text{if } a > b \end{cases} \quad (6)$$

The value type of  $\delta$  is three-value logic  $\mathbb{M} = \{T, F, 0\}$ .

It is useful to define a shorthand:

$$\delta f(a \rightarrow b) = \delta(f(a) \rightarrow f(b)) \quad (7)$$

For  $x \in \mathbb{N}$ , define

$$f'(x) = \delta f(x \rightarrow x + 1) \quad (8)$$

For  $x \in \mathbb{B}$ , define (one of the most elegant contributions of the paper)

$$f'(x) = \mathbf{xnor}(\delta(x \rightarrow \neg x), \delta f(x \rightarrow \neg x)) \quad (9)$$

If any of the inputs of XNOR are 0, then XNOR evaluates to 0. The interpretation of Equation 9 can be understood as the following example. Consider  $f'(F) = T$ . The truth on the right-hand-side indicates the value of the function  $f(x)$  evaluated at  $x = F$  will vary to the same direction as  $x$  varies. I.e. when  $x$  varies from  $F$  to  $T$  ( $x$  is increasing), the value of  $f(x)$  will also increase. That is to say,  $f(T) > f(F)$ .

Here is another example, consider  $f'(T) = F$ . This indicates the variant of the function value will go to the opposite direction as  $x$  varies. Thus as the value of  $x$  going down (i.e.  $T \rightarrow F$ ), the value of  $f(x)$  will go up. I.e.  $f(F) > f(T)$ .

Here is another example, consider  $f'(T) = 0$ . This indicates changing  $x$  from T to F will not affect value of  $f(x)$ , i.e.  $f(T) = f(F)$ .

### 2.2. Commonly use results.

The following useful shorthand will be very helpful while conducting chain rule.

$$f'(x) = \frac{\delta f(x)}{\delta x} \quad (10)$$

State without proof the following properties:

$$\forall x \in \mathbb{B}, \frac{\delta \mathbf{xnor}(a, x)}{\delta x} = a \quad (11)$$

$$\forall x \in \mathbb{B}, \frac{\delta \mathbf{xor}(a, x)}{\delta x} = \neg a \quad (12)$$

The chain rule is calculated (another one of the most elegant contributions of the paper) as:

$$\begin{aligned} \forall x \in \mathbb{B}, (g \circ f)'(x) &= g(f(x))' \\ &= \mathbf{xnor}(g'(f(x)), f'(x)) \end{aligned} \quad (13)$$

Eq. Equation 13 answers the questions: how will the loss change if a boolean weight is flipped.

Lastly, state without proof, the derivative on the cutoff activation function Equation 3

$$\sigma'_\tau(x) = \frac{\delta \sigma_\tau(x)}{\delta x} = \begin{cases} T & \text{if } x = \tau - 1 \\ 0 & \text{else} \end{cases} \quad (14)$$

Notice, the equation above answers how the function value will change if  $x$  is increased by 1. This is because  $\sigma$  is a natural-number-input function, not a boolean-input function. It only increases as  $x$  increases from  $\tau - 1$  to  $\tau$ . Now we are ready to present a complete example.

### 3. A CONCRETE EXAMPLE

First state the configuration. Each sample is a three-dimensional boolean vector, i.e.  $x_k^1 \in \mathbb{B}^3$ . As an example,  $x_{k,1}^1 \in \mathbb{B}$  is first input of the  $k^{\text{th}}$  sample in the batch (remember we index from 1). It is also the first input of the first layer of the  $k^{\text{th}}$  sample. The first layer is  $3 \times 2$ . The first layer not only do the weight multiplication, it applies the activation functions as well. The second layer is  $2 \times 1$ . The second layer's output is the final prediction. This prediction will compare with the  $k^{\text{th}}$  ground truth label  $y_k$ . Recall the loss for the  $k^{\text{th}}$  sample is  $\text{loss}_k = \mathbf{xor}(y_k, x_{k,1}^2)$ .

#### 3.1. Forward Path (FP).

Only consider the  $k^{\text{th}}$  sample in the batch. The first layer operation is (1) applying weight multiplication, and (2) applying the activation. Recall  $L$  is  $\mathbf{xnor}$ . The weight multiplication is the following:

$$\begin{aligned} s_{k,1}^1 &= w_{0,1}^1 + L(w_{1,1}^1, x_{k,1}^1) + \\ &\quad L(w_{2,1}^1, x_{k,2}^1) + L(w_{3,1}^1, x_{k,3}^1) \end{aligned} \quad (15)$$

$$\begin{aligned} s_{k,2}^1 &= w_{0,2}^1 + L(w_{1,2}^1, x_{k,1}^1) + \\ &\quad L(w_{2,2}^1, x_{k,2}^1) + L(w_{3,2}^1, x_{k,3}^1) \end{aligned} \quad (16)$$

Now calculate the activation to get the next layer input:

$$x_{k,1}^2 = \sigma_\tau^1(s_{k,1}^1) \quad (17)$$

$$x_{k,2}^2 = \sigma_\tau^1(s_{k,2}^1) \quad (18)$$

Recall that  $\tau$  is the number of input for the input layer, thus in this case we have  $\tau = \lceil \frac{3}{2} \rceil = 2$  for  $\sigma_\tau^1$ . Now we are ready to calculate the second layer.

The weight multiplication is:

$$s_{k,1}^2 = w_{0,1}^2 + L(w_{1,1}^2, x_{k,1}^2) + L(w_{2,1}^2, x_{k,2}^2) \quad (19)$$

The activation is:

$$x_{k,1}^3 = \sigma_\tau^2(s_{k,1}^2) \quad (20)$$

$x_{k,1}^3$  is also the model prediction. Finally, we can calculate the sample-wise loss

$$\text{loss}_k = \mathbf{xor}(y_k, x_{k,1}^3) \quad (21)$$

The total loss is:

$$\mathcal{L} = \sum_{k=1}^K \text{loss}_k(y_k, x_{k,1}^3) \quad (22)$$

Where  $K$  is the total number of samples in a batch.

### 3.2. Backward Path (BP).

We first consider the BP for a single sample, then aggregate to the batch. Now assume we want to answer the question: if  $w_{1,1}^2$  is flipped, will  $\text{loss}_k$  increase or decrease? First, calculate its BP.

$$\begin{aligned} \frac{\delta \text{loss}_k}{\delta w_{1,1}^2} &= \frac{\delta}{\delta w_{1,1}^2} \left( \text{loss}_k(y_k, x_{k,1}^3) \right) \\ &= \mathbf{xnor} \left( \frac{\delta \text{loss}_k(y_k, x_{k,1}^3)}{\delta x_{k,1}^3}, \frac{\delta x_{k,1}^3}{\delta w_{1,1}^2} \right) \\ &= \mathbf{xnor} \left( \neg y_k, \frac{\delta x_{k,1}^3}{\delta w_{1,1}^2} \right) \\ &= \mathbf{xnor} \left( \neg y_k, \mathbf{xnor} \left( \frac{\delta x_{k,1}^3}{\delta s_{k,1}^2}, \frac{\delta s_{k,1}^2}{\delta w_{1,1}^2} \right) \right) \\ &= \mathbf{xnor} \left( \neg y_k, \mathbf{xnor} \left( \sigma_\tau^{2'}(s_{k,1}^2), x_{k,1}^2 \right) \right) \end{aligned} \quad (23)$$

The expression above can be evaluated already, because  $s_{k,1}^2$  and  $x_{k,1}^2$  have been calculated in FP.

Now let's calculate the variant on weights in the first layer.

$$\begin{aligned}
\frac{\delta \text{loss}_k}{\delta w_{3,2}^1} &= \mathbf{xnor} \left( \frac{\delta \text{loss}_k(y_k, x_{k,1}^3)}{\delta x_{k,1}^3}, \frac{\delta x_{k,1}^3}{\delta w_{3,2}^1} \right) \\
&= \mathbf{xnor} \left( \neg y_k, \mathbf{xnor} \left( \frac{\delta x_{k,1}^3}{\delta s_{k,1}^2}, \frac{\delta s_{k,1}^2}{\delta w_{3,2}^1} \right) \right) \\
&= \mathbf{xnor} \left( \neg y_k, \mathbf{xnor} \left( \sigma_{\tau}^{\prime 2}(s_{k,1}^2), \mathbf{xnor} \left( \frac{\delta s_{k,1}^2}{\delta x_{k,2}^2}, \frac{\delta x_{k,2}^2}{\delta w_{2,3}^1} \right) \right) \right)
\end{aligned} \tag{24}$$

The equation is awfully long. Now only computer the component

$$\begin{aligned}
\mathbf{xnor} \left( \frac{\delta s_{k,1}^2}{\delta x_{k,2}^2}, \frac{\delta x_{k,2}^2}{\delta w_{2,3}^1} \right) &= \mathbf{xnor} \left( w_{2,1}^2, \mathbf{xnor} \left( \frac{\delta x_{k,2}^2}{\delta s_{k,2}^1}, \frac{\delta s_{k,2}^1}{\delta w_{2,3}^1} \right) \right) \\
&= \mathbf{xnor} \left( w_{2,1}^2, \mathbf{xnor} \left( \sigma_{\tau}^{\prime 1}(s_{k,2}^1), x_{k,3}^1 \right) \right)
\end{aligned} \tag{25}$$

At this point, the value can be fully evaluated.

Let's go through a final example to illustrate how to handle the variant of summations. Consider a three-layer network.  $1 \times 2, 2 \times 2, 2 \times 1$ . For simplicity, we drop the subscript indicating the sample index. The final output  $x_1^4$ . The loss is  $x_1^4 \oplus y$  ( $\oplus$  is the notation for XOR). When conducting backward propagation, we run into this term

$$s_1^3 = w_0^3 + \mathbf{xnor}(w_{1,1}^3, x_1^3) + \mathbf{xnor}(w_{2,1}^3, x_2^3) \tag{26}$$

This express has one key difference compared with the previous examples, because if we want to calculate  $\frac{\delta \text{loss}}{\delta w_{1,2}^1}$ , the expression will have two terms remaining. In all previous examples, only one term will remain. Explicitly

$$\frac{\delta s_1^3}{\delta w_{1,2}^1} = \frac{\delta \mathbf{xnor}(w_{1,1}^3, x_1^3)}{\delta w_{1,2}^1} + \frac{\delta \mathbf{xnor}(w_{2,1}^3, x_2^3)}{\delta w_{1,2}^1} \tag{27}$$

in which both  $x_1^3$  and  $x_2^3$  depend on  $w_{1,2}^1$ . Appendix Proposition A.2 of the paper [1] indicates the summation can be conducted by type casting the three-value logic into integer. The rule is:

$$e(a) = \begin{cases} +1 & \text{if } a = T \\ 0 & \text{if } a = 0 \\ -1 & \text{if } a = F \end{cases} \tag{28}$$

### 3.3. Optimization.

Finally we are ready to update the weights on a per-batch basis. Let's say we want to decide whether  $w_{1,2}^1$  should be flipped. The procedure is the following: first calculate

$$q_{1,2,k}^1 = \frac{\delta \text{loss}_k}{\delta w_{1,2}^1} \tag{29}$$

Now we aggregate each sample in the batch. For the  $k^{\text{th}}$  sample, if  $q_{1,2,k}^1 = T$ , this indicates the loss will vary to the same direction as  $w_{1,2}^1$  varies. 0 means it will not vary as the weight varies. False indicates it will vary the opposite direction. Thus

we need to count the number of True and False in the batch, so that we can tell the net result of flipping the weight. The paper has an elegant formulation on this. First we calculate the aggregated result

$$q_{1,2}^1 = (\text{num. of True samples}) - (\text{num. of False samples}) \quad (30)$$

Then the net decision is simply

$$w_{1,2}^1 = \neg w_{1,2}^1 \text{ if } \mathbf{xnor}(q_{1,2}^1, w_{1,2}^1) = T \quad (31)$$

Now, the BP can be conducted fully (hopefully). The paper has more contents regarding how to handle integer mix type and some other optimizers.

### 3.4. Boolean Function Differentiation.

Let's take another approach to calculus of variations, but this time purely based on booleans and boolean functions. Instead of using three-valued logic to indicate direction of change, we can use a single boolean to indicate whether a function's output would change when its input is flipped.

#### 3.4.1. Definition.

For any boolean function  $f : \mathbb{B} \rightarrow \mathbb{B}$ , we define its boolean derivative at point  $x$  as:

$$\frac{\delta f}{\delta x} = \text{xor}(f(x), f(\neg x)) \quad (32)$$

This derivative has a simple interpretation:

- If  $\frac{\delta f}{\delta x} = T$ , then flipping  $x$  will change the output of  $f$
- If  $\frac{\delta f}{\delta x} = F$ , then flipping  $x$  will not change the output of  $f$

#### 3.4.2. Properties.

The boolean derivative has several useful properties:

1. **Symmetry:** The derivative is symmetric around the flip point

$$\frac{\delta f}{\delta x} = \frac{\delta f}{\delta \neg x} \quad (33)$$

2. **Chain Rule:** For composed functions  $g \circ f$ , the derivative follows:

$$\frac{\delta g \circ f}{\delta x} = \text{AND}\left(\frac{\delta f}{\delta x}, \frac{\delta g}{\delta f(x)}\right) \quad (34)$$

To verify this, consider  $f(x) = \text{not } x$  and  $g(x) = \text{and}(T, x)$ . Then:

$$(g \circ f)(x) = \text{and}(T, \text{not } x) = \text{not } x \quad (35)$$

We can verify the chain rule:

- Left side:  $\frac{\delta g \circ f}{\delta x} = \text{xor}(\text{not } x, \text{not}(\text{not } x)) = T$
- Right side:  $\text{xnor}\left(\frac{\delta f}{\delta x}, \frac{\delta g}{\delta f(x)}\right) = \text{xnor}\left(T, \frac{\delta g}{\delta \text{not } x}\right) = \text{xnor}\left(T, \text{xor}(\text{and}(T, \text{not } x), \text{and}(T, \text{not}(\text{not } x)))\right) = \text{xnor}(T, \text{xor}(\text{not } x, x)) = \text{xnor}(T, T) = T$

### 3. Basic Function Derivatives:

- For constant functions:  $\frac{\delta c}{\delta x} = F$
- For identity function:  $\frac{\delta \text{id}}{\delta x} = T$
- For NOT function:  $\frac{\delta \text{not}}{\delta x} = T$
- For AND function:  $\frac{\delta \text{and}(a, x)}{\delta x} = a$

- For OR function:  $\frac{\delta \text{or}(a,x)}{\delta x} = \neg a$
- For XOR function:  $\frac{\delta \text{xor}(a,x)}{\delta x} = T$
- For XNOR function:  $\frac{\delta \text{xnor}(a,x)}{\delta x} = T$

### 3.4.3. Simple Example.

Let's consider a simple boolean circuit with two weights  $w_1, w_2$  and one input  $x$ :

$$f(x, w_1, w_2) = \text{or}(w_1, \text{and}(w_2, x)) \quad (36)$$

To find how the output would change if we flip each weight, we calculate:

1. For  $w_1$  (with fixed  $w_2$  and  $x$ ):

$$\frac{\delta f}{\delta w_1} = \text{xor}(f(w_1, w_2, x), f(\neg w_1, w_2, x)) \quad (37)$$

Let's evaluate this for  $x = T, w_2 = T$ :

$$\frac{\delta f}{\delta w_1} = \text{xor}(\text{or}(w_1, \text{and}(T, T)), \text{or}(\neg w_1, \text{and}(T, T))) \quad (38)$$

$$= \text{xor}(\text{or}(w_1, T), \text{or}(\neg w_1, T)) \quad (39)$$

$$= \text{xor}(T, T) = F \quad (40)$$

This makes sense - when  $x = T$  and  $w_2 = T$ , the output is always T regardless of  $w_1$ .

2. For  $w_2$  (with fixed  $w_1$  and  $x$ ):

$$\frac{\delta f}{\delta w_2} = \text{xor}(f(w_1, w_2, x), f(w_1, \neg w_2, x)) \quad (41)$$

Let's evaluate this for  $x = T, w_1 = F$ :

$$\frac{\delta f}{\delta w_2} = \text{xor}(\text{or}(F, \text{and}(w_2, T)), \text{or}(F, \text{and}(\neg w_2, T))) \quad (42)$$

$$= \text{xor}(w_2, \neg w_2) = T \quad (43)$$

This makes sense - when  $x = T$  and  $w_1 = F$ , flipping  $w_2$  will always change the output.

Let's verify our results using the chain rule:

1. For  $w_1$  with  $x = T, w_2 = T$ :

$$\frac{\delta f}{\delta w_1} = \frac{\delta \text{or}(w_1, \text{and}(w_2, x))}{\delta w_1} \quad (44)$$

Since this is a direct OR with  $w_1$ , we can compute:

$$= \frac{\delta \text{or}(w_1, T)}{\delta w_1} \quad (45)$$

$$= \begin{cases} -T & \text{if } w_1 = F \\ T & \text{if } w_1 = T \end{cases} \quad (46)$$

$$= F \quad (47)$$

2. For  $w_2$  with  $x = T, w_1 = F$ : Here we need the chain rule since  $w_2$  appears inside the AND:

$$\frac{\delta f}{\delta w_2} = \text{and}\left(\frac{\delta \text{and}(w_2, x)}{\delta w_2}, \frac{\delta \text{or}(w_1, \_)}{\delta \text{and}(w_2, x)}\right) \quad (48)$$

With  $x = T, w_1 = F$ :

$$= \text{and}\left(\frac{\delta \text{and}(w_2, T)}{\delta w_2}, \frac{\delta \text{or}(F, \_)}{\delta \text{and}(w_2, T)}\right) \quad (49)$$

Let's compute each part:

- First part:  $\frac{\delta \text{and}(w_2, T)}{\delta w_2} = \begin{cases} T & \text{if } w_2 = F \\ -T & \text{if } w_2 = T \end{cases} = T$
- Second part:  $\frac{\delta \text{or}(F, y)}{\delta y} \big|_{y=\text{and}(w_2, T)} = \begin{cases} -F & \text{if } y = F \\ F & \text{if } y = T \end{cases} = T$

Therefore:

$$= \text{and}(T, T) = T \quad (50)$$

Both methods give us the same results, confirming our calculations.

#### REFERENCES

1. Nguyen, V. M., Ocampo, C., Askri, A., Tran, B.-H.: Boolean Logic for Low-Energy Deep Learning. In: 2nd Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ICML 2024) (2024)

DEPARTMENT OF COMPUTER SCIENCE, BOSTON UNIVERSITY, BOSTON, U.S.A.  
*Email address:* wfchen@bu.edu