

---

# On the Inductive Bias of Neural Networks for Learning Read-once DNFs (Supplementary material)

---

Ido Bronstein<sup>1</sup>

Alon Brutzkus<sup>1</sup>

Amir Globerson<sup>1</sup>

<sup>1</sup>Blavatnik School of Computer Science , Tel Aviv University, Israel

Here we provide complete proofs for the results in the main paper, as well as additional empirical results.

## 1 GRADIENT FLOW - ASSUMPTIONS

When gradient flow is implemented on non-differentiable functions (e.g., ReLU) the implementation can choose from among a set of possible sub-differentials. Here we define which of these our analysis will use. This choice corresponds to the common way gradient methods are implemented in practice for the ReLU function.

Recall that the gradient flow step is  $\frac{d\theta^{(t)}}{dt} \in -\partial^\circ L(\theta^{(t)})$  for a.e.  $t$ , where:

$$\partial^\circ f(\mathbf{x}) = \text{conv} \left\{ \lim_{k \rightarrow \infty} \nabla f(\mathbf{x}_k) \mid \mathbf{x}_k \rightarrow \mathbf{x} \text{ and } f \text{ is differentiable at } \mathbf{x}_k \right\} \quad (7)$$

is the Clarke's subdifferential.

As discussed in the main text, we will assume that the gradient flow step selects a specific vector in the subdifferential. This is done by setting the subgradient of ReLU at 0 in advance to a constant value  $a \in [0, 1]$ . Namely, this value of the subgradient is used for all neurons and in all iterations. Usually  $a$  is set to be either 0 or 1.

Formally, for each  $i \in r$  we denote  $\frac{d\mathbf{w}_i^{(t)}}{dt} = \frac{1}{m} \sum_{\mathbf{x} \in \mathbb{S}} \frac{d\mathbf{w}_i^{(t)}(\mathbf{x})}{dt}$ . Here,  $\frac{d\mathbf{w}_i^{(t)}(\mathbf{x})}{dt}$  is the gradient update of  $\mathbf{w}_i^{(t)}$  restricted to the summand that depends on  $\mathbf{x}$  in  $L$  (Eq. (2) of the main paper). Similarly, we denote  $\frac{db_i^{(t)}}{dt} = \frac{1}{m} \sum_{\mathbf{x} \in \mathbb{S}} \frac{db_i^{(t)}(\mathbf{x})}{dt}$ . For our result, we need the following technical assumption.

**Assumption 1.1.** *There exists an  $a \in [0, 1]$  such that for every step  $t > 0$ , every neuron  $i \in [r]$ , and every sample  $\mathbf{x} \in \mathbb{S}$  if  $\mathbf{w}_i^{(t)} \cdot \mathbf{x} + b_i^{(t)} = 0$  then  $\frac{d\mathbf{w}_i^{(t)}(\mathbf{x})}{dt} = ay\ell' (yN(\mathbf{x}; \theta^{(t)})) \mathbf{x}$  and  $\frac{db_i^{(t)}}{dt} = ay\ell' (yN(\mathbf{x}; \theta^{(t)}))$ .*

## 2 PROOF OF THEOREM 5.1

We first need the following notation and recall the KKT conditions in our context. Let  $\partial^\circ \sigma(\mathbf{w}_i \cdot \mathbf{x}_l + b) \subseteq \mathbb{R}^{D+1}$  be the subdifferential of neuron  $i$  given input  $\mathbf{x}_l$ . It holds that:

$$\partial^\circ \sigma(\mathbf{w}_i \cdot \mathbf{x}_l + b) = \begin{cases} \{(\mathbf{x}_l, 1)\} & \text{if } \mathbf{w}_i \cdot \mathbf{x}_l + b_i > 0 \\ \{(\mathbf{0}, 0)\} & \text{if } \mathbf{w}_i \cdot \mathbf{x}_l + b_i < 0 \\ [0, 1]^{D+1} & \text{if } \mathbf{w}_i \cdot \mathbf{x}_l + b_i = 0 \end{cases}$$

**KKT conditions:** A feasible point  $\theta = (\mathbf{W}, \mathbf{b}, c)$  of the min norm problem (Eq. (4) of the main paper) is a KKT point if there exist  $\lambda_1, \dots, \lambda_m \geq 0$  such that:

1. **Stationarity:**

$$\forall i \in [r], \mathbf{w}_i = \sum_{l \in [m]} \lambda_l y_l \mathbf{h}_{il} \text{ and } b_i = \sum_{l \in [m]} \lambda_l y_l g_{il} \quad (8)$$

and

$$c = \sum_{j \in [m]} \lambda_j y_j \quad (9)$$

where  $(\mathbf{h}_{il}, g_{il}) \in \partial^\circ \sigma(\mathbf{w}_i \cdot \mathbf{x}_l + b)$ .

2. **Complementary slackness:** if  $y_i N(\mathbf{x}_i, \theta) > 1$ , then  $\lambda_i = 0$ .

We now proceed to prove the theorem. In the first part, we show that by Assumption 1.1, we can restrict the possible values of  $\mathbf{h}_{il}, g_{il}$  for a KKT point that GF converges to. In the second part, we prove properties of neurons that memorize samples. In the third part, we use the previous parts to show that memorizing solutions cannot be KKT points.

**Part 1:** In this part we prove the following lemma.

**Lemma 2.1.** *Assume that Assumption 1.1 holds and GF converges to a KKT point with parameters  $\mathbf{h}_{il}$  and  $g_{il}$  for  $1 \leq i \leq r$  and  $1 \leq l \leq n$ , as defined in Eq. (8). Then,  $(\mathbf{h}_{il}, g_{il}) \in \{(\mathbf{x}_l, 1), (a\mathbf{x}_l, a), (0, 0)\}$ .*

*Proof.* For each  $1 \leq l \leq n, 1 \leq i \leq r$  and  $t > 0$  let  $(\mathbf{h}_{il}^{(t)}, g_{il}^{(t)}) \in \partial^\circ \sigma(\mathbf{w}_i^{(t)} \cdot \mathbf{x}_l + b_i^{(t)})$  be the corresponding values of the GF step at time  $t$  in the subdifferential  $\partial^\circ \sigma(\mathbf{w}_i^{(t)} \cdot \mathbf{x}_l + b_i^{(t)})$ . By inspecting the proofs of Lyu & Li (2020) and Dutta et al. (2013), we see that  $(\mathbf{h}_{il}, g_{il})$  is equal to the limit of a convergent subsequence of  $\left\{ \left( \mathbf{h}_{il}^{(t_j)}, g_{il}^{(t_j)} \right) \right\}_{j=0}^\infty$ .

By assumption 1.1, we know that for each  $j$ ,  $(\mathbf{h}_{il}^{(t_j)}, g_{il}^{(t_j)}) \in \{(\mathbf{x}_l, 1), (a\mathbf{x}_l, a), (0, 0)\}$ . Therefore, the limit also satisfies  $(\mathbf{h}_{il}, g_{il}) \in \{(\mathbf{x}_l, 1), (a\mathbf{x}_l, a), (0, 0)\}$ . □

**Part 2:** We will first need the following definition.

**Definition 2.1.** *Given a sample  $\hat{\mathbf{x}} \in \mathcal{X}$  and index  $j \in [D]$ , the sample with Hamming distance one from  $\hat{\mathbf{x}}$  at index  $j$  is defined as  $\mathcal{H}(\hat{\mathbf{x}}, j) \in \mathcal{X}$  and satisfies the following:*

$$\mathcal{H}(\hat{\mathbf{x}}, j)_j = -x_j \text{ and } \forall j' \in [D] \setminus \{j\} \mathcal{H}(\hat{\mathbf{x}}, j)_{j'} = x_{j'} \quad (10)$$

*The set of all samples with Hamming distance one from  $\hat{\mathbf{x}}$  is defined as:  $\Psi(\hat{\mathbf{x}}) = \{\mathbf{x}' \in \mathcal{X} \mid \exists j \in [D] \mathcal{H}(\hat{\mathbf{x}}, j) = \mathbf{x}'\}$ . Note,  $|\Psi(\hat{\mathbf{x}})| = D$*

Using this definition, we rephrase Lemma 5.1 of the main text and prove those properties of memorizing neurons (Definition 5.1 of the main paper):

**Lemma 2.2.** *Let  $D > 2$ . If a neuron  $i \in [r]$  memorizes a sample  $\hat{\mathbf{x}} \in \mathbb{S}_x$ , then it satisfies the following properties:*

1.  $\hat{x}_j = \text{sign}(w_{ij})$  for all  $1 \leq j \leq D$ .
2. For  $\mathbf{x} \in \mathcal{X}$  if  $\mathbf{w}_i \cdot \mathbf{x} + b_i = 0$  then  $\mathbf{x} \in \Psi(\hat{\mathbf{x}})$ .
3.  $b_i < 0$

*Proof. Property 1:* Assume by contradiction that there exists  $j \in [D]$  such that  $\text{sign}(w_{ij}) \neq \hat{x}_j$ . Then  $\mathbf{w}_i \cdot \mathcal{H}(\hat{\mathbf{x}}, j) + b_i \geq \mathbf{w}_i \cdot \hat{\mathbf{x}} + b_i > 0$ , in contradiction to the memorization assumption in Eq. (5) of the main paper.

**Property 2:** Assume by contradiction that there exists  $\mathbf{x} \in \mathcal{X} \setminus (\Psi(\hat{\mathbf{x}}) \cup \{\hat{\mathbf{x}}\})$  such that  $\mathbf{w}_i \cdot \mathbf{x} + b_i = 0$ . We define  $J = \{j \in [D] \mid x_j = -\hat{x}_j\}$  and  $j' \in J$  for some index in  $J$ . By Property 1, the sample  $\tilde{\mathbf{x}} = \mathcal{H}(\hat{\mathbf{x}}, j')$  satisfies the

following:

$$\begin{aligned}
\mathbf{w}_i \cdot \mathbf{x} - \mathbf{w}_i \cdot \tilde{\mathbf{x}} &= \sum_{j \in [D] \setminus J} w_{ij} x_j + \sum_{j \in J} w_{ij} x_j - \sum_{j \in [D] \setminus \{j^0\}} w_{ij} \tilde{x}_j - w_{ij^0} \tilde{x}_j \\
&= \sum_{j \in [D] \setminus J} w_{ij} \hat{x}_j - \sum_{j \in J} w_{ij} \hat{x}_j - \sum_{j \in [D] \setminus \{j^0\}} w_{ij} \hat{x}_j + w_{ij^0} \hat{x}_j \\
&= \sum_{j \in [D] \setminus J} |w_{ij}| - \sum_{j \in J} |w_{ij}| - \sum_{j \in [D] \setminus \{j^0\}} |w_{ij}| + |w_{ij^0}| = -2 \sum_{j \in J \setminus \{j^0\}} |w_{ij}|
\end{aligned} \tag{11}$$

Since  $\mathbf{x} \notin \Psi(\hat{\mathbf{x}})$ , it holds that  $J \setminus \{j^0\} \neq \emptyset$ . Furthermore, by Property 1,  $\forall j \in J \setminus \{j^0\}$   $w_{ij} \neq 0$ . Thus,  $\sum_{j \in J \setminus \{j^0\}} |w_{ij}| > 0$  which implies that  $\mathbf{w}_i \cdot \mathbf{x} < \mathbf{w}_i \cdot \tilde{\mathbf{x}}$ . We know that  $\mathbf{w}_i \cdot \mathbf{x} + b_i = 0$ , therefore  $\mathbf{w}_i \cdot \tilde{\mathbf{x}} + b_i > 0$ . This contradicts the memorization assumption in Eq. (5) of the main paper.

**Property 3:** Assume by contradiction that  $b_i \geq 0$ . We define  $j' = \arg \min_{j \in [D]} \{|w_{ij}|\}$ . For the sample  $\tilde{\mathbf{x}} = \mathcal{H}(\hat{\mathbf{x}}, j')$ , the following holds by Property 1:

$$\mathbf{w}_i \cdot \tilde{\mathbf{x}} = \sum_{j \in [D] \setminus \{j^0\}} w_{ij} \hat{x}_j - w_{ij^0} \hat{x}_{j^0} = \sum_{j \in [D] \setminus \{j^0\}} |w_{ij}| - |w_{ij^0}| > 0 \tag{12}$$

Since  $D > 2$ , we have  $\sum_{j \in [D] \setminus \{j^0\}} |w_{ij}| - |w_{ij^0}| > 0$ . Thus,  $\mathbf{w}_i \cdot \tilde{\mathbf{x}} + b_i > 0$  which contradicts the memorization assumption in Eq. (5) of the main paper.  $\square$

**Part 3:** Now, we can proceed to prove Theorem 5.1.

Consider a network with parameters  $\bar{\boldsymbol{\theta}} = (\bar{\mathbf{W}}, \bar{\mathbf{b}}, \bar{c})$ , neuron  $i \in [r]$  and a sample  $(\hat{\mathbf{x}}, \hat{y}) \in \mathbb{S}$  such that Eq. (5) of the main paper holds (i.e., the neuron  $i$  memorizes the sample  $\hat{\mathbf{x}}$ ). We assume by contradiction that there exists an initialization  $\boldsymbol{\theta}^{(0)}$  such that if we run gradient flow from  $\boldsymbol{\theta}^{(0)}$  using  $\mu$  then the weights  $\boldsymbol{\theta}^{(t)}$  will converge to  $\bar{\boldsymbol{\theta}}$ . According to the results of Lyu & Li (2020); Ji & Telgarsky (2020), we know that there exists  $\alpha > 0$  such that  $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{b}, c) = \alpha \bar{\boldsymbol{\theta}}$  is a KKT point of Eq. (4) of the main paper and Eq. (8) holds for  $\boldsymbol{\theta}$ . Note that for  $\boldsymbol{\theta}$ , neuron  $i$  memorizes the sample  $(\hat{\mathbf{x}}, \hat{y})$  as well.

Given a sample  $\mathbf{x}_l \in \mathbb{S}_x$  we can see that the following holds:

1. If  $\mathbf{x}_l = \hat{\mathbf{x}}$  then  $(\mathbf{h}_{il}, g_{il}) = \{(\hat{\mathbf{x}}, 1)\}$ .
2. If  $\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})$  then  $(\mathbf{h}_{il}, g_{il}) = \{(a\mathbf{x}_l, a)\}$  where  $a \in [0, 1]$  by Assumption 1.1.
3. If  $\mathbf{x}_l \notin \Psi(\hat{\mathbf{x}})$  then  $(\mathbf{h}_{il}, g_{il}) = \{(\mathbf{0}, 0)\}$  by Property 2 of Lemma 2.2.

We will show that for every  $\lambda_1, \dots, \lambda_{|\mathbb{S}|} \geq 0$ , Eq. (8) does not hold. We can assume without loss of generality that all the samples in  $\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x$  are support vectors and  $(\mathbf{h}_{il}, g_{il}) = (\mathbf{x}_l, 1)$ . This is because if one of the samples is not a support vector then we can take  $\lambda_l = 0$ . Furthermore, if  $(\mathbf{h}_{il}, g_{il}) = (\mathbf{0}, 0)$  then we can take  $\lambda_l = 0$  and if  $(\mathbf{h}_{il}, g_{il}) = (a\mathbf{x}_l, a)$  for  $a > 0$  we can set  $\frac{\lambda_l}{a}$  instead of  $\lambda_l$ . Under this assumption we can write Eq. (8) for  $\bar{\boldsymbol{\theta}}$  using only  $\lambda_1, \dots, \lambda_D$  and  $\hat{\lambda}$  that correspond to the samples of  $\Psi(\hat{\mathbf{x}})$  and  $\hat{\mathbf{x}}$ , respectively:

$$\mathbf{w}_i = \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l \mathbf{x}_l + \hat{\lambda} \hat{y} \hat{\mathbf{x}} \quad \text{and} \quad b_i = \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l + \hat{\lambda} \hat{y} \tag{13}$$

Assume by contradiction that there exists  $\mathbf{x}_{l^0} = \mathcal{H}(\hat{\mathbf{x}}, j')$  such that  $\mathbf{x}_{l^0} \in \Psi(\hat{\mathbf{x}}) \setminus \mathbb{S}_x$ . The following holds by Eq. (13) and Definition 2.1:

$$\begin{aligned}
(1) \quad w_{ij^0} &= \hat{\lambda} \hat{y} \hat{x}_{j^0} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l \hat{x}_{j^0} \\
(2) \quad b_i &= \hat{\lambda} \hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l
\end{aligned} \tag{14}$$

Using the first property in Lemma 2.2:

$$\begin{aligned}
(1) \quad w_{ij^0} &= \hat{\lambda}\hat{y} \operatorname{sign}(w_{ij^0}) + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l \operatorname{sign}(w_{ij^0}) \\
(2) \quad b_i &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l
\end{aligned} \tag{15}$$

Therefore,

$$\begin{aligned}
(1) \quad |w_{ij^0}| &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l \\
(2) \quad b_i &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \cap \mathbb{S}_x} \lambda_l y_l
\end{aligned} \tag{16}$$

This means that  $0 \leq |w_{ij^0}| = b_i$ , which is in contradiction to the third property of Lemma 2.2. Therefore, we can assume from now on that  $\Psi(\hat{\mathbf{x}}) \subseteq \mathbb{S}_x$ , and we can write the KKT conditions as follows:

$$\mathbf{w}_i = \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})} \lambda_l y_l \mathbf{x}_l + \hat{\lambda}\hat{y}\hat{\mathbf{x}} \quad \text{and} \quad b_i = \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})} \lambda_l y_l + \hat{\lambda}\hat{y} \tag{17}$$

Given  $\mathbf{x}_{l^0} = \mathcal{H}(\hat{\mathbf{x}}, j')$ , the following holds by Eq. (17) and Definition 2.1:

$$\begin{aligned}
(1) \quad w_{ij^0} &= \hat{\lambda}\hat{y}\hat{x}_{j^0} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \setminus \{\mathbf{x}_{l^0}\}} \lambda_l y_l \hat{x}_{j^0} - \lambda_{l^0} y_{l^0} \hat{x}_{j^0} \\
(2) \quad b_i &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})} \lambda_l y_l
\end{aligned} \tag{18}$$

Using the first property in Lemma 2.2:

$$\begin{aligned}
(1) \quad w_{ij^0} &= \hat{\lambda}\hat{y} \operatorname{sign}(w_{ij^0}) + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \setminus \{\mathbf{x}_{l^0}\}} \lambda_l y_l \operatorname{sign}(w_{ij^0}) - \lambda_{l^0} y_{l^0} \operatorname{sign}(w_{ij^0}) \\
(2) \quad b_i &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})} \lambda_l y_l
\end{aligned} \tag{19}$$

Therefore,

$$\begin{aligned}
(1) \quad |w_{ij^0}| &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \setminus \{\mathbf{x}_{l^0}\}} \lambda_l y_l - \lambda_{l^0} y_{l^0} \\
(2) \quad b_i &= \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}})} \lambda_l y_l
\end{aligned} \tag{20}$$

The result of subtracting (2) – (1) is:

$$b_i - |w_{ij^0}| = 2\lambda_{l^0} y_{l^0} \tag{21}$$

Next we show that it must hold that  $y_{l^0} = -1$ . To see this, note that the third property of Lemma 2.2 implies  $b_i < 0$  and therefore  $b_i - |w_{ij^0}| < 0$ . Assuming in contradiction that  $y_{l^0} = 1$ , the RHS of Eq. (21) satisfies  $0 \leq 2\lambda_{l^0} y_{l^0}$ . But we just saw that the LHS satisfies  $b_i - |w_{ij^0}| < 0$ . We therefore have a contradiction and conclude that  $y_{l^0} = -1$ . We can conclude that  $\Psi(\hat{\mathbf{x}})$  contains only negative samples.

Next we argue that  $\hat{y} = -1$ . To see this, assume in contradiction that  $\hat{y} = 1$ . Then there exists a  $n \in [K]$  such that  $\hat{\mathbf{x}}$  satisfies the term  $t_n^*$ . Due to the fact that  $K \geq 2$ , we know that there exists  $j \in [D] \setminus \mathbb{A}_n$ . Then,  $\mathcal{H}(\hat{\mathbf{x}}, j) \in \Psi(\hat{\mathbf{x}})$  is a positive sample in contradiction to the fact that  $\Psi(\hat{\mathbf{x}})$  contains only negative samples, and therefore  $\hat{y} = -1$ .

By Eq. (21) we know that for every  $j \in [D]$  a sample  $\mathbf{x}_l = \mathcal{H}(\hat{\mathbf{x}}, j)$  satisfies the following:  $\lambda_l = \frac{1}{2}(|w_{ij}| - b_i)$ . If we assign this in Eq. (20), we get:

$$|w_{ij^0}| = \hat{\lambda}\hat{y} + \sum_{\mathbf{x}_l \in \Psi(\hat{\mathbf{x}}) \setminus \{\mathbf{x}_{l^0}\}} \frac{1}{2}(b_i - |w_{ij}|) - \frac{1}{2}(b_i - |w_{ij^0}|) \tag{22}$$

Therefore,

$$0 = \hat{y} + \frac{1}{2}(D - 2)b - \frac{1}{2} \sum_j w_{ij} |j|_1 \quad (23)$$

But this results in a contradiction, because we know that  $0, \frac{1}{2}(D - 2)b < 0$  and  $\frac{1}{2} \sum_j w_{ij} |j|_1 < 0$ .

Thus, we conclude that gradient flow cannot converge to memorizing solutions.

### 3 PROOF OF THEOREM 6.1

We prove the theorem in several parts. We first prove properties of a perfect solution (Section 3.1). In Section 3.2 we prove several results regarding the bias threshold. In Section 3.3 we prove auxiliary lemmas and in Section 3.4 we prove the alignment of the neurons of the optimal solutions to the terms of the DNF. We conclude the proof in Section 3.5.

#### 3.1 A SIMPLE PROPERTY OF PERFECT SOLUTIONS

Recall the definitions  $S_+ = \{x \in \{0, 1\}^n : \sum_j x_j \geq 2\}$  and  $S_- = \{x \in \{0, 1\}^n : \sum_j x_j \leq 1\}$ . We first need the following definitions. We say that a solution  $(W; b)$  satisfies the  $\text{MIN}_+$  property if for any positive point  $x \in S_+$  there exists  $\delta \in [r]$  such that  $w_\delta x + b_\delta \geq 2$ . We say that a solution satisfies the  $\text{MIN}_-$  property if for any negative point  $x \in S_-$  and for all  $i \in [r]$ ,  $w_i x + b_i \leq 0$ .

Lemma 3.1.  $(W; b)$  is a perfect solution if and only if  $(W; b)$  satisfies  $\text{MIN}_+$  and  $\text{MIN}_-$ .

Proof. If  $(W; b)$  is a perfect solution, then for all  $(k; y) \in S_+$ ,  $y_N(x; W; b) \geq 1$ . Therefore, if  $y = 1$ ,  $\sum_{i \in [r]} (w_i x + b_i) \geq 2$ . Thus, there exists  $\delta \in [r]$  such that  $w_\delta x + b_\delta \geq 2$  and the  $\text{MIN}_+$  holds. If  $y = 0$ , then  $\sum_{i \in [r]} (w_i x + b_i) \leq 0$  and therefore for all  $i \in [r]$ ,  $w_i x + b_i \leq 0$ . The other direction follows similarly.  $\square$

We note that one direct consequence of Lemma 3.1 is that given a negative  $x \in S_-$  is activated by a neuron  $i \in [r]$ , i.e.,  $w_i x + b_i > 0$ , then the  $\text{MIN}_-$  property doesn't hold, and is not a perfect solution.

#### 3.2 PROOF OF LEMMA 6.1

In this section we show that when  $S_n = X$ , the bias of any neuron in a perfect solution is upper bounded by a certain value which we call the bias threshold. To simplify the formulation of this section we define the following:

Definition 3.1. We define the set of indices which are not active in any term as the noisy indices and denote them by  $A_{K+1} = [D] \setminus \bigcup_{n \in [K]} A_n$ .

Definition 3.2. For each term  $n \in [K]$  of  $f$  and  $i \in [r]$  define  $V_n(w_i) = \max_{j \in A_n} \min_{j \in A_n} w_{ij} g_j$ .

Definition 3.3. The bias threshold for a weight is  $\text{BT}(w) = \sum_j |w_{jj}|_1 + 2 \sum_{n \in [K]} V_n(w)$ .

Note,  $\text{BT}(w) \geq 0$  because every term includes at least 2 literals. Using those definitions we rephrase Lemma 6.1 and prove it:

Lemma 3.2. Assume that  $S_n = X$ .  $(W; b)$  satisfies that  $\forall i \in [r] \quad b_i \leq \text{BT}(w_i)$  and satisfies the  $\text{MIN}_+$  property if and only if the network is a perfect solution.

Proof. We will show that given a neuron  $(w; b)$ , there is a negative sample  $x \in X$  for which  $w x + b > 0$  if and only if  $b > \text{BT}(w)$ . By showing that and using the assumption  $S_n = X$ , we can conclude that  $\forall i \in [r] \quad b_i \leq \text{BT}(w_i)$  if and only if the  $\text{MIN}_-$  property holds. By combining this with Lemma 3.1, we can prove our claim.

Given  $\text{neuron}(w; b)$ , we define the minimum index of a term  $\in [K]$  as  $J_n = \arg \min_{j \in [K]} w_j$ . Consider a sample  $x \in \mathbb{R}^D$  that is defined by:

$$x_j = \begin{cases} \text{sign}(w_j) & \text{if } V_n(w) > 0 \wedge j = J_n \\ \text{sign}(w_j) & \text{otherwise} \end{cases} \quad (24)$$

For every term  $\in [K]$ , if  $V_n(w) > 0$  then  $x_{J_n} = \text{sign}(w_{J_n}) = -1$ . Otherwise,  $V_n(w) = 0$  and there exists  $\exists j \in [K]$  such that  $w_j < 0$ , i.e.,  $x_j = \text{sign}(w_j) = -1$ . In any case,  $x_{J_n} < 0$ . Therefore, the label of this sample is negative and denote it by  $y = -1$ .

We show that  $x = \text{BT}(w)$  by,

$$\begin{aligned} w \cdot x &= \sum_{j \in [D]} w_j x_j = \sum_{j \in [K+1]} w_j x_j \\ &= \sum_{j \in [K] \text{ and } V_n(w) > 0} w_j \text{sign}(w_j) + w_{J_n} \text{sign}(w_{J_n}) \\ &+ \sum_{j \in [K] \text{ and } V_n(w) = 0} w_j \text{sign}(w_j) + \sum_{j \in [K+1]} w_j \text{sign}(w_j) \\ &= \sum_{j \in [K] \text{ and } V_n(w) > 0} w_j x_j + \sum_{j \in [K] \text{ and } V_n(w) = 0} w_j x_j + \sum_{j \in [K+1]} w_j x_j \\ &= \sum_{j \in [K] \text{ and } V_n(w) > 0} w_j x_j + 2V_n(w) + \sum_{j \in [K] \text{ and } V_n(w) = 0} w_j x_j + 2V_n(w) + \sum_{j \in [K+1]} w_j x_j \\ &= \sum_{j \in [K]} w_j x_j + 2V_n(w) + \sum_{j \in [K+1]} w_j x_j = \sum_{j \in [D]} w_j x_j = 2 \sum_{j \in [K]} w_j x_j \\ &= 2 \sum_{j \in [K]} w_j x_j = 2 \sum_{j \in [K]} V_n(w) = \text{BT}(w) \end{aligned} \quad (25)$$

For the first direction, if  $b > \text{BT}(w)$  then  $w \cdot x + b = \text{BT}(w) + b > 0$ , as desired.

In the second direction, assume that there is a negative sample  $x$  such that  $w \cdot x + b > 0$ . We will show that  $x = w \cdot x$ . Every term  $\in [K]$  satisfies for all  $j \in [K+1]$ :

$$x_j w_j - |w_j| = \text{sign}(w_j) w_j = x_j w_j \quad (26)$$

If  $V_n(w) = 0$ , then the index  $j = J_n$  also satisfies Eq. (26), by the definition of  $J_n$ . Otherwise,  $V_n(w) > 0$  and we know that there exists  $\exists j^0 \in [K+1]$  such that  $x_{j^0} = -1$  (since  $x$  is negative), and  $w_{j^0} > 0$ . If  $J_n = j^0$ , then  $x_{j^0} w_{j^0} = x_{J_n} w_{J_n}$ . Otherwise, the following holds:

$$x_{j^0} w_{j^0} + x_{J_n} w_{J_n} = w_{j^0} + w_{J_n} < 0 = w_{j^0} - w_{J_n} = x_{j^0} w_{j^0} - x_{J_n} w_{J_n} \quad (27)$$

Note that every index  $j \in [K+1]$  satisfies Eq. (26) as well. Therefore,  $w \cdot x = w \cdot x$ . We can conclude that:

$$0 < w \cdot x + b = w \cdot x + b = \text{BT}(w) + b \quad (28)$$

which implies that  $b > \text{BT}(w)$  as desired.  $\square$

From this point, we will assume  $x = w \cdot x$  without mentioning it explicitly.

### 3.3 AUXILIARY LEMMAS

We first define a special positive sample for every term. The special sample is a sample where all indices corresponding to the term will have positive values, and all other indexes will have negative values. The special samples will be used as the hardest positive samples to satisfy the MIN property.

Definition 3.4. For a term  $n \in [K]$ , we define the special sample  $x^{(n)} \in S_x$  of this term as follows:

$$\forall j \in A_n, x_j^{(n)} = 1 \text{ and } \forall j \in [D] \setminus A_n, x_j^{(n)} = -1 \quad (29)$$

We denote the set of all the special samples by  $S_x = \{x^{(n)} \mid n \in [K]\}$ .

Lemma 3.3. Given  $\mathcal{N} = (W; b)$ , assume the following conditions are satisfied:

1.  $\forall i \in [r]; \forall j \in [D], w_{ij} \geq 0$ .
2. For every  $x \in \mathcal{O}$  there exists  $\delta \in [r]$  such that  $\sum_{i \in I} w_i x + b_i \geq 2$ .

Then  $\mathcal{N}$  satisfies the MIN property.

Proof. Let  $x \in S_x$ . Then  $\exists n \in [K]$  such that  $\forall j \in A_n, x_j = 1$ . By the second assumption  $\exists \delta \in [r]$  such that  $\sum_{i \in I} w_i x^{(n)} + b_i \geq 2$ .

For every  $i \in [r]$  the following holds:

$$w_i x = \sum_{j \in [D]} w_{ij} x_j = \sum_{j \in A_n} w_{ij} + \sum_{j \in [D] \setminus A_n} x_j w_{ij} \quad (30)$$

From the first condition of the claim we can deduce that

$$\sum_{j \in A_n} w_{ij} + \sum_{j \in [D] \setminus A_n} x_j w_{ij} = \sum_{j \in A_n} w_{ij} + \sum_{j \in [D] \setminus A_n} w_{ij} = \sum_{j \in [D]} w_{ij} x_j^{(n)} = w_i x^{(n)} = w_i \quad (31)$$

Then:

$$\sum_{i \in I} (w_i x + b_i) = \sum_{i \in I} (w_i x^{(n)} + b_i) \geq 2 \quad (32)$$

and  $\mathcal{N}$  satisfies the MIN property for  $x$  as required.  $\square$

The following definition will be very useful in our analysis.

Definition 3.5. Given a min-norm solution  $\mathcal{N} = (W; b)$ , we say that a solution  $\hat{\mathcal{N}} = (\hat{W}; \hat{b})$  is a modified solution if the following holds:

$$\forall i \in [r], \forall j \in [D], \hat{w}_{ij} = w_{ij} \text{ and } \hat{b}_i = b_i \quad (33)$$

Thus, given a min-norm solution, to define a modified solution, we only need to define the new  $(\hat{w}_i; \hat{b}_i)$ .

Lemma 3.4. Given a min-norm solution  $\mathcal{N} = (W; b)$ , every  $i \in [r]$  satisfies:

1.  $b_i = BT(w_i)$ .
2.  $\forall j \in [D], w_{ij} \geq 0$ .
3.  $\exists n \in [K]$  such that  $\forall j \in A_n, w_{ij} \geq 0$  and  $\forall j \in [D] \setminus A_n, w_{ij} = 0$ .

Proof. Property 1: Assume by contradiction that  $b_i \notin BT(w_i)$ . By Lemma 3.2,  $b_i$  has to be smaller than  $BT(w_i)$ , because otherwise  $\mathcal{N}$  is not a perfect solution. Now consider the modified solution,  $\hat{\mathcal{N}}$ , which is defined by:

$$\hat{w}_i = w_i; \hat{b}_i = BT(\hat{w}_i) \quad (34)$$

By the assumption  $\hat{b}_i > b_i$ . Then, every  $x \in S_x$  satisfies the following:

$$\sum_i x_i w_i + b_i < \sum_i x_i \hat{w}_i + \hat{b}_i \quad (35)$$

Since  $\hat{x}$  satisfies the MIN<sub>+</sub> property, the above implies that  $\hat{x}$  satisfies it as well.

From Definition 3.5 every  $x \in S$  satisfies:

$$\sum_{i \in [r]} x_i w_{i0} + b_{i0} = \sum_i x_i \hat{w}_{i0} + \hat{b}_{i0} \quad (36)$$

In addition, we saw in the proof of Lemma 3.2 that if  $\hat{x} = BT(\hat{w})$  then  $\sum_i x_i \hat{w}_i + \hat{b}_i = 0$ . Therefore,  $\hat{x}$  satisfies the MIN property. By Lemma 3.1,  $\hat{x}$  is a perfect solution.

From Definition 3.3 the bias threshold is nonpositive and therefore  $\hat{b}_i \leq 0$  implies that  $\hat{w}_i < j b_j$ . We know that  $\hat{w}_i = w_i$  and therefore  $\|(\hat{w}_i; \hat{b}_i)\|_2^2 < \|j w_i; b_j\|_2^2$  which contradicts the optimality of  $\hat{x}$ .

Property 2: Assume by contradiction that  $\exists j \in [D]$  such that  $w_{ij0} < 0$ . Consider the following  $j$ -modified solution  $\hat{x}$ :

$$\hat{x}_j = w_{ij} \wedge \hat{w}_{ij0} = 0 \wedge \hat{b}_i = BT(\hat{w}_i) \quad (37)$$

We want to show that:

$$\sum_{n \in [K]} V_n(w_i) = \sum_{n \in [K]} V_n(\hat{w}_i) \quad (38)$$

If  $\exists n \in [K]$  such that  $j \in A_{n0}$ , then it follows that  $V_{n0}(w_i) = V_{n0}(\hat{w}_i) = 0$  and Eq. (38) is satisfied. Otherwise,  $j \in A_{K+1}$ , by Definition 3.1, the indices of  $A_{K+1}$  don't affect the value of the sums in Eq. (38). Therefore, this equation is satisfied in this case as well.

We know that  $b_i = BT(w_i)$  according to Property 1 above, thus every  $x \in S_x$  satisfies the following:

$$\begin{aligned} \sum_i x_i w_i + b_i &= \sum_j x_j w_{ij} + BT(w_i) \\ &= \sum_{j \in [D]} x_j w_{ij} + \sum_{j \in [D]} x_j w_{ij0} + \sum_{j \in [D]} x_j w_{ij} + \sum_{n \in [K]} 2V_n(w_i) \end{aligned} \quad (39)$$

We can see that  $\sum_{j \in [D]} x_j w_{ij0} + \sum_{j \in [D]} x_j w_{ij} \leq 0$ . Then,

$$\begin{aligned} &\sum_{j \in [D]} x_j w_{ij} + \sum_{j \in [D]} x_j w_{ij0} + \sum_{j \in [D]} x_j w_{ij} + \sum_{n \in [K]} 2V_n(w_i) \\ &= \sum_{j \in [D]} x_j w_{ij} + \sum_{j \in [D]} x_j w_{ij} + \sum_{n \in [K]} 2V_n(w_i) \\ &= \sum_{j \in [D]} x_j \hat{w}_{ij} + \sum_{n \in [K]} 2V_n(\hat{w}_i) = \sum_i \hat{w}_i + BT(\hat{w}_i) = \sum_i \hat{w}_i + \hat{b}_i \end{aligned} \quad (40)$$

Using the fact that  $\sum_i w_i + b_i \leq \sum_i \hat{w}_i + \hat{b}_i$  with the fact that  $\hat{x}$  satisfies the MIN<sub>+</sub> property, we can conclude that  $\hat{x}$  satisfies this property too.

According to Property 1 above and Definition 3.5 we have:

$$\sum_{i \in [r]} BT(\hat{w}_{i0}) = BT(w_{i0}) = b_{i0} = \hat{b}_{i0} \quad (41)$$

In addition, we know that  $\hat{x} = BT(\hat{w}_i)$  by Eq. (37). According to Lemma 3.2, we know that  $\hat{x}$  is a perfect solution.

From Eq. (37), we know that  $\sum_{j \in [D]} x_j w_{ij} > \sum_{j \in [D]} x_j w_{ij0} + \sum_{j \in [D]} x_j w_{ij}$ . Combining this with Eq. (38) and the fact that the bias threshold is nonpositive we can conclude that

$$\sum_{n \in [K]} V_n(w_i) < \sum_{n \in [K]} V_n(\hat{w}_i) + \sum_{j \in [D]} x_j BT(w_{ij}) > \sum_{j \in [D]} x_j BT(\hat{w}_{ij}) + \sum_{j \in [D]} x_j b_j > \hat{b}_i \quad (42)$$



Therefore  $\|j(w_i; \hat{b}_i)\|_2^2 < \|j(w_i; b_i)\|_2^2$  in contradiction to the optimality of  $\hat{b}_i$ .

Property 3: Assume by contradiction that there exists  $\epsilon [r]$  such that:

$$\exists n_1 \in n_2 \subseteq [K + 1] \text{ such that } \exists j \in A_{n_1} w_{ij} > 0 \text{ and } \exists j \in A_{n_2} w_{ij} > 0 \quad (43)$$

Without loss of generality we assume:

$$\sum_{j \in A_{n_1}} w_{ij} > 0 \quad \sum_{j \in A_{n_2}} w_{ij} > 0 \quad (44)$$

Let's look on the following  $j$ -modified  $\hat{b}_i$  which is defined by:

$$\exists j \in A_{n_2} \hat{w}_{ij} = 0 \text{ and } \exists j \in [D] \setminus A_{n_2} \hat{w}_{ij} = w_{ij} \text{ and } \hat{b}_i = BT(\hat{w}_i) \quad (45)$$

First, we will show that  $\hat{b}_i \leq b_i$ . If  $n_2 \in [K + 1]$ , from Definition 3.3 and the assumption that  $n_2 \geq 1$ , the following holds:

$$\hat{b}_i = BT(\hat{w}_i) = BT(w_i) + \sum_{j \in A_{n_2}} j w_{ij} j - 2V_{n_2}(w_i) = b_i + \sum_{j \in A_{n_2}} w_{ij} - 2V_{n_2}(w_i) \leq b_i \quad (46)$$

Otherwise  $n_2 = K + 1$  and from Definition 3.3 the following holds:

$$\hat{b}_i = BT(\hat{w}_i) = BT(w_i) + \sum_{j \in A_{n_2}} j w_{ij} j = b_i + \sum_{j \in A_{n_2}} w_{ij} \leq b_i \quad (47)$$

In both cases  $\hat{b}_i \leq b_i$  as required.

Given  $\epsilon \in [K]$ , we know that  $\hat{b}_i$  is a perfect solution and thus it satisfies the  $\epsilon$ -MIN property. Then, for  $x^{(\epsilon)}$  there exists  $\epsilon [r]$  such that:

$$\sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} \leq \epsilon \quad (48)$$

We will show that there exists  $\epsilon [r]$  such that:

$$\sum_{i \in I} \hat{w}_{i0} x^{(\epsilon)} + \hat{b}_{i0} \leq \epsilon \quad (49)$$

Recall, by Definition 3.5, for any  $\epsilon \in [r]$  we know that  $\sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} = \sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0}$ .

If  $\epsilon \in n_2$ , due to Property 1 and Property 2 above and the fact that  $\hat{b}_i \leq b_i$  the following holds:

$$\begin{aligned} \sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} &= \sum_{j \in [D] \setminus A_{n_2}} w_{ij} x_j^{(\epsilon)} + \sum_{j \in A_{n_2}} w_{ij} + b_{i0} < \sum_{j \in [D] \setminus A_{n_2}} w_{ij} x_j^{(\epsilon)} + b_{i0} + \sum_{j \in [D]} w_{ij} x_j^{(\epsilon)} + \hat{b}_{i0} \\ &= \sum_{i \in I} \hat{w}_{i0} x^{(\epsilon)} + \hat{b}_{i0} \end{aligned} \quad (50)$$

Therefore,

$$\sum_{i \in I} \hat{w}_{i0} x^{(\epsilon)} + \hat{b}_{i0} \leq \sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} \leq \epsilon \quad (51)$$

Otherwise  $\epsilon \in n_2$ . By the fact that  $\hat{b}_i = BT(\hat{w}_i) \leq 0$ , Property 2 above and Eq. (44) the following holds:

$$\sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} = \sum_{j \in A_{n_2}} w_{ij} + \sum_{j \in [D] \setminus A_{n_2}} w_{ij} + b_{i0} = \sum_{j \in A_{n_2}} w_{ij} + \sum_{j \in A_{n_1}} w_{ij} + b_{i0} = \sum_{j \in A_{n_2}} w_{ij} + \sum_{j \in A_{n_1}} w_{ij} \leq 0 \quad (52)$$

Therefore, using Definition 3.5,  $\hat{b}_i$  satisfies the following:

$$\sum_{i \in I} \hat{w}_{i0} x^{(\epsilon)} + \hat{b}_{i0} = \sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} \leq \sum_{i \in I} w_{i0} x^{(\epsilon)} + b_{i0} \leq \epsilon \quad (53)$$

We can conclude that  $\hat{b}_i$  satisfies Eq. (49). Combining this with Property 2 above, we can see that  $\hat{b}_i$  meets the condition of Lemma 3.3 and then it satisfies the MIN property.

According to Property 1 above and Definition 3.5:

$$\inf_{i \in [r]} \text{BT}(w_{i0}) = \text{BT}(w_{i0}) = b_0 = \hat{b}_0 \quad (54)$$

In addition, we know that  $\hat{b}_0 = \text{BT}(w_i)$  by Eq. (45). According to Lemma 3.2, the solution  $\hat{b}_0$  is a perfect solution.

As we saw  $0 \leq \hat{b}_0 - b_0 \leq \hat{b}_0 - b_j, \forall j \in [r], \exists j \in [r] \text{ such that } w_{ij} = 0 = w_{ij}$  and  $\exists j \in [r] \text{ such that } w_{ij} > 0 = w_{ij}$  and therefore  $\|j(w_i; b)\|_2^2 > \|j(w_i; \hat{b}_0)\|_2^2$  which contradicts the optimality of  $\hat{b}_0$ .  $\square$

### 3.4 ALIGNMENT LEMMAS

The following three lemmas show the alignment properties of the min-norm solution.

Lemma 3.5. Given a min-norm solution  $(w; b)$ , every  $i \in [r]$  either aligns with some term  $j \in [K]$  or it holds that  $w_i = 0; b_j = 0$

Proof. Given  $i \in [r]$ , if  $w_i = 0$  the claim is true by Property 1 of Lemma 3.4. Otherwise, by Property 3 of Lemma 3.4,  $\exists j \in [K]$  such that:

$$\exists j \in [r] \text{ such that } w_{ij} = 0 \text{ and } \exists j \in [r] \text{ such that } w_{ij} > 0 \text{ and } \exists j \in [r] \text{ such that } w_{ij} = 0 \quad (55)$$

Assume by contradiction that:

$$\exists j_1, j_2 \in [r] \text{ such that } w_{ij_1} \in w_{ij_2} \quad (56)$$

Without loss of generality, we assume  $w_{ij_1} > w_{ij_2}$  and  $w_{ij_2} = \min_{j \in [r]} w_{ij}$ .

Define the following  $i$ -modified solution  $\hat{w}$ :

$$\hat{w}_{ij} = w_{ij_2} \text{ and } \hat{w}_{ij} = w_{ij} \text{ and } \hat{b}_0 = \text{BT}(\hat{w}_i) \quad (57)$$

Note that  $V_n(w_i) = V_n(\hat{w}_i) = w_{ij_2}$ .

Given  $x \in S_+$ , we know that  $(w; b)$  satisfies the MIN property. Then  $\exists i \in [r]$  such that:

$$\sum_{i \in [r]} w_{i0} x + b_0 \leq 2 \quad (58)$$

Recall, by Definition 3.5, for any  $i \in [r]$  we know that  $w_{i0} x + \hat{b}_0 = w_{i0} x + b_0$

If  $x \cdot t_n = \|j t_n\|_1$ , due to Property 3 of Lemma 3.4 the following holds:

$$\begin{aligned} w_i x + b_0 &= \sum_{j \in [r]} w_{ij} + \text{BT}(w_i) = \sum_{j \in [r]} w_{ij} + 2V_n(w_i) = 2w_{ij_2} \\ &= \sum_{j \in [r]} \hat{w}_{ij} + 2w_{ij_2} = \sum_{j \in [r]} \hat{w}_{ij} + \|j t_n\|_1 + 2V_n(\hat{w}_i) = \hat{w}_i x + \hat{b}_0 \end{aligned} \quad (59)$$

Then we can conclude:

$$\sum_{i \in [r]} \hat{w}_{i0} x + \hat{b}_0 = \sum_{i \in [r]} w_{i0} x + b_0 \leq 2 \quad (60)$$

Otherwise  $x \cdot t_n < \|j t_n\|_1$  and by Definition 3.3:

$$w_i x + b_0 = \sum_{j \in [r]} w_{ij} x_j + \text{BT}(w_i) = \sum_{j \in [r] \setminus j_2} w_{ij} + \sum_{j \in [r] \setminus j_2} w_{ij} + 2w_{ij_2} = 0 \quad (61)$$

Therefore,

$$\sum_{i \in [r]} \hat{w}_{i0} x + \hat{b}_0 = \sum_{i \in [r]} w_{i0} x + b_0 \leq \sum_{i \in [r]} w_{i0} x + b_0 \leq 2 \quad (62)$$

We can conclude that  $\hat{w}_i$  satisfies the  $\text{MIN}_+$  property.

According to Property 1 of Lemma 3.4 and Definition 3.5:

$$\|w_{i_0}\|_2 \leq \|w_{i_0}\|_2 \text{ and } \|b_0\|_2 \leq \|b_0\|_2 \quad (63)$$

In addition, we know that  $\hat{w}_i = \text{BT}(w_i)$  by Eq. (57). According to Lemma 3.2, the solution  $\hat{w}_i$  is a perfect solution.

Finally, we can see that  $\|w_{i_1}\|_2 > \|w_{i_2}\|_2$  implies that  $\|\text{BT}(w_{i_1})\|_2 > \|\text{BT}(w_{i_2})\|_2$  and  $\|b_{i_1}\|_2 > \|b_{i_2}\|_2$ . Thus, we have  $\|w_{i_1}\|_2 > \|w_{i_2}\|_2$ . This is in contradiction to the optimality of  $\hat{w}_i$ , as desired.

We can now define  $w_i = V_n(w_i)$  and we know that the neurons satisfy:

$$\|w_{i_1}\|_2 \leq \|w_{i_2}\|_2 \text{ and } \|b_{i_1}\|_2 \leq \|b_{i_2}\|_2 \quad (64)$$

Therefore,  $w_i = w_{i_1} t_n; b_i = b_{i_1} (2 - \|t_n\|_2)$  and we can say that neuron  $i$  aligns the term.  $\square$

Lemma 3.6. Given a min-norm solution  $w = (W; b)$ , every 2 neurons  $i_1, i_2 \in [r]$  that align with term 2  $[K]$  satisfy  $\|w_{i_1}\|_2 = \|w_{i_2}\|_2$ .

Proof. Given  $i_1, i_2 \in [r]$  that align with term 2  $[K]$  we have  $w_{i_1} = w_{i_1} t_n$  and  $w_{i_2} = w_{i_2} t_n$ . Assume by contradiction that  $\|w_{i_2}\|_2 < \|w_{i_1}\|_2$ . Define the following solution  $\hat{w}$ :

$$\begin{aligned} \|w_{i_1}\|_2 &= \|w_{i_1}\|_2; \hat{b}_i = b_i \\ w_{i_1} &= \frac{\|w_{i_1}\|_2 + \|w_{i_2}\|_2}{2} t_n; \hat{b}_{i_1} = \text{BT}(w_{i_1}) \\ w_{i_2} &= \frac{\|w_{i_1}\|_2 + \|w_{i_2}\|_2}{2} t_n; \hat{b}_{i_2} = \text{BT}(w_{i_2}) \end{aligned} \quad (65)$$

Note,  $\|w_{i_0}\|_2 \leq \|w_{i_0}\|_2 + \|b_0\|_2 = \|w_{i_0}\|_2 + \|b_0\|_2$

Given  $x \in S_+$ , we know that  $\hat{w}$  satisfies the  $\text{MIN}_+$  property. Then  $\exists \theta \in [r]$  such that:

$$\|w_{i_0}\|_2 \leq \|w_{i_0}\|_2 + \|b_0\|_2 \quad (66)$$

If  $\|x\|_2 = \|t_n\|_2$  we can calculate the following:

$$\|w_{i_1}\|_2 \|x\|_2 + \|b_{i_1}\|_2 \|x\|_2 + \|w_{i_2}\|_2 \|x\|_2 + \|b_{i_2}\|_2 \|x\|_2 = 2 \|w_{i_1}\|_2 \|x\|_2 + \|b_{i_1}\|_2 \|x\|_2 + \|w_{i_2}\|_2 \|x\|_2 + \|b_{i_2}\|_2 \|x\|_2 \quad (67)$$

Therefore,

$$\|w_{i_0}\|_2 \|x\|_2 + \|b_0\|_2 \|x\|_2 = \|w_{i_0}\|_2 \|x\|_2 + \|b_0\|_2 \|x\|_2 \quad (68)$$

Otherwise  $\|x\|_2 < \|t_n\|_2$  thus  $\|w_{i_1}\|_2 \|x\|_2 + \|b_{i_1}\|_2 \|x\|_2 > 0$ ,  $\|w_{i_2}\|_2 \|x\|_2 + \|b_{i_2}\|_2 \|x\|_2 < 0$ . Then,

$$\|w_{i_0}\|_2 \|x\|_2 + \|b_0\|_2 \|x\|_2 = \|w_{i_0}\|_2 \|x\|_2 + \|b_0\|_2 \|x\|_2 + \|w_{i_0}\|_2 \|x\|_2 + \|b_0\|_2 \|x\|_2 \quad (69)$$

Combining this with Property 2 of Lemma 3.4 we can see that  $\hat{w}$  meets the condition of Lemma 3.3 and then it satisfies the  $\text{MIN}_+$  property.

According to Property 1 of Lemma 3.4:

$$\|w_{i_0}\|_2 \leq \|w_{i_0}\|_2 \text{ and } \|b_0\|_2 \leq \|b_0\|_2 \quad (70)$$

In addition, we know that  $\hat{w}_{i_1} = \text{BT}(w_{i_1})$  and  $\hat{w}_{i_2} = \text{BT}(w_{i_2})$  by Eq. (65). According to Lemma 3.2, the solution  $\hat{w}$  is a perfect solution.

<sup>1</sup>  $w_i$  were defined in the proof of the previous lemma.

We will prove that  $\sum_{i \in I_2} \|w_i; b_i\|_2^2 > \sum_{i \in I_2} \|w_i; \hat{b}_i\|_2^2$ . This will contradict the optimality of  $(w_i; b_i)$ . Note that  $\|w_i; b_i\|_2^2 = \frac{1}{i} |A_n| + \frac{1}{i} (|A_n| - 2)^2$ . Then:

$$\begin{aligned} \sum_{i \in I_2} \|w_i; b_i\|_2^2 &= \sum_{i \in I_2} \|w_i; \hat{b}_i\|_2^2 + \sum_{i \in I_2} \|w_i; b_i\|_2^2 - \sum_{i \in I_2} \|w_i; \hat{b}_i\|_2^2 \\ &= \sum_{i \in I_2} \left( \frac{1}{i} + \frac{1}{i} - 2 \frac{i_1 + i_2}{2} \right) |A_n| + (|A_n| - 2)^2 \\ &= \frac{1}{2} \sum_{i \in I_2} \left( \frac{1}{i_1} - \frac{1}{i_2} \right) |A_n| + (|A_n| - 2)^2 \\ &= \frac{1}{2} \sum_{i \in I_2} \left( \frac{1}{i_1} - \frac{1}{i_2} \right) |A_n| + (|A_n| - 2)^2 \end{aligned}$$

For  $i_1 \in I_2$ , we get  $\sum_{i \in I_2} \|w_i; b_i\|_2^2 - \sum_{i \in I_2} \|w_i; \hat{b}_i\|_2^2 > 0$ , as needed.  $\square$

**Lemma 3.7.** Given a min-norm solution  $(w; b)$ , if  $I_2$  is the set of all neurons that align with term 2 [K], then  $\sum_{i \in I_2} w_i = 0$ .

**Proof.** Assume by contradiction that  $\sum_{i \in I_2} w_i \neq 0$ . If  $\sum_{i \in I_2} w_i < 0$ , then  $\sum_{i \in I_2} w_i \neq 0$ , by Lemma 3.5 we know that the neuron  $i_0$  aligns with another term or is equal to 0, then  $w_{i_0} = 0$ . By Property 1 of Lemma 3.4,  $b_{i_0} = BT(w_{i_0}) = 0$ . Therefore  $x^{(n)} w_{i_0} + b_{i_0} = 0$ . Then, for every  $i \in I_2$  the following holds:

$$\sum_{i \in I_2} w_i x^{(n)} + b_i = \sum_{i \in I_2} w_i x^{(n)} + b_i = \sum_{i \in I_2} |A_n| i + (2 - |A_n|) i = 2 \sum_{i \in I_2} i < 2 \quad (71)$$

and thus  $(w; b)$  doesn't satisfy the MIN property. By Lemma 3.1, this contradicts the fact that  $(w; b)$  is a perfect solution.

If  $\sum_{i \in I_2} w_i > 0$ , we choose an arbitrary  $i_1 \in I_2$ . Define  $\hat{w}$  as follows:

$$\begin{aligned} \forall i \in I_2 \setminus \{i_1\} \quad \hat{w}_i &= w_i, \quad \hat{b}_i = b_i \\ \forall i \in I_2 \setminus \{i_1\} \quad \hat{w}_i &= 0, \quad \hat{b}_i = 0 \\ \hat{w}_{i_1} &= t_n, \quad \hat{b}_{i_1} = BT(\hat{w}_{i_1}) \end{aligned} \quad (72)$$

Given  $x^{(e)} \geq 0$ , we know that  $(\hat{w}; \hat{b})$  satisfies the MIN property. Thus,  $\exists i_0 \in I_2$  such that:

$$\sum_{i \in I_2} \hat{w}_i x^{(e)} + \hat{b}_i \leq 2 \quad (73)$$

We will show that there exists  $i_0 \in I_2$  such that:

$$\sum_{i \in I_2} \hat{w}_i x^{(e)} + \hat{b}_i \leq 2 \quad (74)$$

If  $n = 1$ , then:

$$\hat{w}_{i_1} x^{(e)} + \hat{b}_{i_1} = t_n x^{(e)} + BT(t_n) = |A_n| t_n + |A_n| t_n = 2 t_n \leq 2 \quad (75)$$

Therefore, by choosing  $i_0 = i_1$ , we can show that  $(\hat{w}; \hat{b})$  satisfies Eq. (74).

Otherwise  $n \geq 2$  and then every  $i \in I_2$  satisfies:

$$w_i x^{(e)} + b_i = |A_n| i - (|A_n| - 2) i < 0 \quad (76)$$

From Eq. (72) we can conclude:

$$\sum_{i \in I_2} \hat{w}_i x^{(e)} + \hat{b}_i = \sum_{i \in I_2} w_i x^{(e)} + b_i + \sum_{i \in I_2} w_i x^{(e)} + b_i \leq 2 \quad (77)$$

Therefore,  $\hat{w}$  satisfies Eq. (74) for  $\kappa^{(e)}$ . In addition, by Property 2 of Lemma 3.4 and Eq. (72), we know that all the weights of the neurons in  $\hat{I}$  are nonnegative. Then  $\hat{w}$  meets the condition of Lemma 3.3 and it satisfies the MLP property.

According to Property 1 of Lemma 3.4 and Eq. (72):

$$\sum_{i \in [r]} |b_i - \text{BT}(w_i)| = \sum_{i \in [r]} |b_i - \hat{b}_i| \quad (78)$$

In addition, we know that  $\sum_{i \in [r]} \hat{b}_i = 0 = \text{BT}(w_i)$  and  $\hat{b}_i = \text{BT}(w_i)$ . Therefore, according to Lemma 3.2, the solution  $\hat{w}$  is a perfect solution.

We will prove that  $\sum_{i \in [r]} \|w_i - b_i\|_2^2 > \sum_{i \in [r]} \|w_i - \hat{w}_i\|_2^2$ . This will contradict the optimality of  $\hat{w}$ . Indeed:

$$\begin{aligned} \sum_{i \in [r]} \|w_i - b_i\|_2^2 - \sum_{i \in [r]} \|w_i - \hat{w}_i\|_2^2 &= \sum_{i \in [r]} \|w_i - b_i\|_2^2 - \sum_{i \in [r]} \|w_i - \hat{w}_i\|_2^2 = \\ &= \sum_{i \in [r]} \left( \|w_i - b_i\|_2^2 - \|w_i - \hat{w}_i\|_2^2 \right) = \\ &= \sum_{i \in [r]} \left( \|w_i - b_i\|_2^2 - \|w_i - \hat{w}_i\|_2^2 \right) = \\ &= \sum_{i \in [r]} \left( \|w_i - b_i\|_2^2 - \|w_i - \hat{w}_i\|_2^2 \right) = \end{aligned} \quad (79)$$

Since  $\sum_{i \in [r]} \alpha_i > 1$ , we have  $\sum_{i \in [r]} \|w_i - b_i\|_2^2 - \sum_{i \in [r]} \|w_i - \hat{w}_i\|_2^2 > 0$ , which completes the proof. □

### 3.5 FINISHING THE PROOF OF THEOREM 6.1

Proof. Given a min-norm solution  $w = (W; b)$ , by Lemma 3.5, each neuron  $i \in [r]$  aligns with some term  $t_i \in [K]$  or it is equal to 0. Assume by contradiction that there exists a term  $t \in [K]$  that is not aligned, namely  $\sum_{i \in [r]} \mathbb{1}_{t_i = t} < n$ . Consider the special positive sample  $x^{(n)} \in S_+$ . From the definition of  $x^{(n)}$ , every  $j \in [D]$  satisfies  $x_j^{(n)} = 1$ . Then,

$$\sum_{i \in [r]} x^{(n)} w_i = \sum_{i \in [r]} x^{(n)} t_{n_i} = \sum_{i \in [r]} \mathbb{1}_{t_i = t} < 0 \quad (80)$$

By the first property of Lemma 3.5,  $\sum_{i \in [r]} b_i = \text{BT}(w_i) = 0$ . Therefore  $\sum_{i \in [r]} x^{(n)} w_i + b_i < 0$ , in contradiction to the fact that  $w$  is a perfect solution.

In conclusion, every term  $t \in [K]$  aligns with a set of neurons  $I_t$ . By Lemma 3.7, we know that  $\sum_{i \in I_t} \alpha_i = 1$ . By Lemma 3.6, we know that every two neurons  $i_1, i_2$  that align with the same term satisfy  $\alpha_{i_1} = \alpha_{i_2}$ . Therefore,  $w$  is a DNF recovery solution. □

## 4 EXPERIMENT DETAILS AND ADDITIONAL RESULTS

Selected read-once DNFs: In this section we provide additional experiments for the following types of read-once DNFs.

- $f_1$  - 3-term read-once DNF where the length of every term is 3 ~~for 9~~
- $f_2$  - 4-term read-once DNF where the length of the terms are 4,5,5,5 ~~to 25~~
- $f_3$  - 8-term read-once DNF where the length of the terms are 3,3,4,4,4,4,5,5 ~~to 40~~
- $f_4$  - 10-term read-once DNF where the length of the terms are 3,3,4,4,4,4,6,6,6,6 ~~to 50~~
- $f_5$  - 15-term read-once DNF where of every term is 5 ~~to 100~~

General details: In all the experiments, “small initialization” refers to initializing weights from  $w^{(0)} \sim \mathcal{N}(0; 10^{-6})$  and  $b^{(0)} = [0]^D$ . The learning rate for SGD is  $\eta = 10^{-3}$ , the number of hidden units is  $n = 2000$  and the batch size is 32. We create the train set by sampling uniformly from  $\{0,1\}^D$ . All the experiments can run on any single GPU. Training a single network can take up to two hours.

(a)

(b)

(c)

(d)

Figure 5: Test accuracy for the convex network with small initialization, convex network with large initialization and standard networks. Figure (a) shows the performance when learning for  $f_3$ , Figure (b) for  $f_3$ , Figure (c) for  $f_4$  and Figure (d) for  $f_3$  (Results for  $f_2$  were presented in the main paper).

Weight Matrix Visualization: When presenting weight matrices, we first cluster the neurons using the Hierarchical clustering algorithm<sup>2</sup>. We then plot the weight values in an image, where neurons clustered together appear in consecutive rows. Note this of course does not change the model itself, but makes it easy to see if there are well clustered neurons (as in the DNF recovery case).

Sample Complexity Experiments: We evaluate test accuracy as a function of the training sample size for different models. Results are shown in Figure 5. Specifically, we compare the convex network with small initialization (see details above), the convex network with large initialization (we take  $W^{(0)} \sim N(0; 1)$ ), and a “standard” network with one hidden layer, same width as the convex network and Xavier initialization (we checked different initialization schemes, including small Gaussian initialization, and verified that this not affect the results). We run every experiment 100 times and present the mean performance and the std of this mean (the std is small and smaller than the line width).

Implementation of the Statistical Queries (SQ) Methods: In Figure 2a of the main paper, we present results of the statistical query algorithm. We implemented the algorithm described in Mansour & Schain (2001). We view  $\epsilon$  as a hyperparameter. Therefore, we use 10% of the train set as validation for  $n$ . We present the performance of the SQ algorithm only for  $D = 9$ , because for larger dimension the algorithm failed in creating a DNF for the range of train set sizes tested.

<sup>2</sup>We used `scipy.cluster.hierarchy.linkage` with `centroid` as a method.

(a) (b) (c) (d)

Figure 6: (a-c) Effect of training size on the learned model (see learned models for different training sizes in the text) for the ground-truth model  $\theta_1$ . Panels a-c correspond to training sizes 800, 1500 and 7500. (d) Result for training on 7500 with large initialization.

(a) (b) (c) (d)

Figure 7: Same as Figure 6 but with  $\theta_2$  as ground truth.

(a) (b) (c) (d)

Figure 8: Same as Figure 6 but with  $\theta_3$  as ground truth.

(a) (b) (c) (d)

Figure 9: Same as Figure 6 but with  $\theta_4$  as ground truth.

---

**Algorithm 1** Reconstruction Procedure

---

**Input:** Network  $\theta = (\mathbf{W}, \mathbf{b}, c)$ , DNF  $f^*$ , fixed sets  $A, B \subseteq [0, 1]^L$ .  
**Output:** True if the network with parameter  $\theta$  reconstructs DNF  $f^*$ , False otherwise.

```
for  $(a, b) \in A \times B$  do
   $\mathbf{W}' \leftarrow [0]^{r \times D}$  ▷  $\mathbf{W}'$  will be a  $\{0, 1\}$  matrix where each row corresponds to a term in a DNF.
  for  $1 \leq i \leq r$  do
    if  $\ell_\infty(\mathbf{w}_i) \geq a * \ell_\infty(\mathbf{W})$  then ▷ Taking into account only meaningful neurons
      for  $1 \leq j \leq D$  do
        if  $w_{ij} \geq \ell_\infty(\mathbf{w}_i) * b$  then ▷ Taking into account only meaningful values
           $w'_{ij} \leftarrow 1$ 
        end if
      end for
    end if
  end for
end for
if The set of terms represented by  $\mathbf{W}'$  is exactly the set of terms of the DNF  $f^*$  then
  return True
end if
end for
return False
```

---

**Learned models for different training sizes:** In the main paper, we show empirically that learning convex networks with small init and GD leads to a DNF recovery solution, and we also show formally that in the population risk DNF recovery is norm minimizing. Here we show explicit model weights for different training sizes, demonstrating that approximate DNF recovery solutions are obtained for fairly small sample sizes. Figures 6,7,8,9,10 (panels a-c) show these results. In panel d of these figures we show the learned model for when learning with large Gaussian initialization and with the same train set size. It can be seen that larger initialization does not result in the recovery-DNF solution (note we are also visualizing these solutions using clustering as explained above, and there is clearly no cluster structure in the solution).

**DNF reconstruction:** In Figure 2b of the main paper we present accuracy results for DNF reconstruction. To obtain these, we take the learned model and apply a simple rounding procedure to check if this model reconstructs the ground-truth DNF. The procedure is outlined in Algorithm 1. In the procedure, we create a  $\{0, 1\}$  matrix  $\mathbf{W}'$  where the column indices of 1s in each row correspond to a term of a DNF. Thus,  $\mathbf{W}'$  represents a set of terms. If the set of terms of  $\mathbf{W}'$  is exactly the set of terms of the input DNF, the procedure returns True. In our experiments we ran the procedure with inputs  $A = [0, 0.1, 0.2, \dots, 0.9]$  and  $B = [0, 0.2, 0.4, \dots, 0.8]$ .

**The effect of learning  $c$ :** Figure 11 shows that fixing the learnable parameter  $c$  to  $-1$  does not effect the structure of the solution. In this experiment, we took 2 networks with the same width: One with learnable  $c$  initialized to 0, and the second with fixed  $c = -1$ . We initialize the other weights with the same values, and train them with the same train set for the same number of steps. Finally we plot the solution that the network learns. We can say they both recover the underlying DNF.

**Tabular datasets:** We consider the three UCI datasets: kr-vs-kp, Splice, and diabetes. For these, we convert the input into binary by changing categorical variables to one-hot. We also consider binary classification such that in kr-vs-kp the class 'won' is positive considered and 'notwon' is negative, in Splice the classes 'EI' and 'IE' are considered positive and 'N' negative, and diabetes is binary by design. We train on 90% of the data and test on 10%. The reconstruction process is identical to algorithm 1 when instead of validate if  $\mathbf{W}'$  is identical to  $f^*$ , we return the  $\mathbf{W}'$  with the best accuracy on the train set.

The relevant code can be found in our repository: <https://github.com/idobronstein/Exploring-the-Inductive-Bias-of-Neural-Networks-for-Learning-Read-once-DNFs.git>.



