

1 APPENDIX

1.1 Derivation of (BETALINEAR) rule:

For dense-based neural networks with *linear* activation function, for a given neural network with a layer $N : y = \text{linear}(\bar{w}.\bar{x} + \bar{b})$ and $\delta : y \bowtie \bar{n}$, the computed precondition δ' is $\bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b}$.

Given, postcondition $\delta : y \bowtie \bar{n}$ and given $N : y = \text{linear}(\bar{w}.\bar{x} + \bar{b})$ We know from the weakest precondition equation $\delta' = \text{wp}(N, \delta)$.

$$\begin{aligned} \text{So, } \delta' &= \text{wp}(N, \delta) \equiv \text{wp}(\text{linear}(\bar{w}.\bar{x} + \bar{b}), y \bowtie \bar{n}) \\ &\equiv \text{linear}(\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \equiv (\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \\ &\equiv \bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b}; [\text{From least square solution}] \\ &\equiv \bar{x} \bowtie (\bar{y}.\bar{n}) - \bar{b}; [\bar{y} = (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T)] \end{aligned} \quad (2.1.1)$$

1.2 Derivation of (BETARELU) rule:

For dense-based neural network with *relu* activation function, for given $N : y = \text{relu}(\bar{w}.\bar{x} + \bar{b})$ and $\delta : y \bowtie \bar{n}$, the computed precondition δ' is $\bar{x} \bowtie ((\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b}) \wedge \bar{x} \bowtie ((\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).(-\bar{b}))$.

Given, postcondition $\delta' : y \bowtie \bar{n}$ and given $N : y = \text{relu}(\bar{w}.\bar{x} + \bar{b})$ We know from the weakest precondition equation $\delta' = \text{wp}(N, \delta)$.

$$\begin{aligned} \text{So, } \delta' &= \text{wp}(N, \delta) \equiv \text{wp}(\text{relu}(\bar{w}.\bar{x} + \bar{b}), y \bowtie \bar{n}) \\ &\equiv \text{relu}(\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \equiv \max(0, \bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \\ \text{We know, if } m &= \max(x, y) \text{ then,} \\ (m = x \wedge (x &\geq y) \vee (m = y) \vee (y \geq x)) \\ \text{So, } \delta' &\equiv \max(0, \bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \\ &\equiv \{0 \wedge 0 \geq (\bar{w}.\bar{x} + \bar{b})\} \vee \{(\bar{w}.\bar{x} + \bar{b}) \wedge (\bar{w}.\bar{x} + \bar{b}) \geq 0\} \bowtie \bar{n} \end{aligned} \quad (2.1.2)$$

$$\begin{aligned} \text{From (2.1.2), } 0 \wedge 0 &\geq (\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \\ &\equiv 0 \leq \bar{n} \wedge 0 \geq (\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \equiv \bar{w}.\bar{x} + \bar{b} \bowtie \bar{n} \\ &\equiv \bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b} \end{aligned} \quad (2.1.3)$$

$$\begin{aligned} \text{Again, } \{(\bar{w}.\bar{x} + \bar{b}) \wedge (\bar{w}.\bar{x} + \bar{b}) &\geq 0\} \bowtie \bar{n} \\ \text{So, } \bar{w}.\bar{x} + \bar{b} \bowtie \bar{n} &\equiv \bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b}; [\text{From (2.1.1)}] \\ \wedge \bar{w}.\bar{x} + \bar{b} \geq 0 &\bowtie \bar{n} \equiv \bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).(-\bar{b}) \end{aligned} \quad (2.1.4)$$

$$\begin{aligned} \text{From (2.1.3) \& (2.1.4), } P &\equiv \bar{x} \bowtie ((\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\bar{n} - \bar{b}) \\ &\wedge \bar{x} \bowtie ((\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).(-\bar{b})) \end{aligned}$$

1.3 Derivation of (BETASIGMOID) rule:

For a dense-based neural network with *sigmoid* activation function, for a given neural network with a single layer $N : y = \text{sigmoid}(\bar{w}.\bar{x} + \bar{b})$ and $\delta : y \bowtie \bar{n}$, The computed precondition δ' is $\bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\ln(\frac{\bar{n}}{1-\bar{n}}) - \bar{b}$.

Given, postcondition $\delta : y \bowtie \bar{n}$ and given $N : y = \text{sigmoid}(\bar{w}.\bar{x} + \bar{b})$ We know from the weakest precondition equation, $\delta' = \text{wp}(N, \delta)$.

$$\begin{aligned} \text{So, } \delta' &= \text{wp}(N, \delta) \equiv \text{wp}(\text{sigmoid}(\bar{w}.\bar{x} + \bar{b}), y \bowtie \bar{n}) \\ &\equiv \text{sigmoid}(\bar{w}.\bar{x} + \bar{b}) \bowtie \bar{n} \equiv \frac{1}{1 + e^{-(\bar{w}.\bar{x} + \bar{b})}} \bowtie \bar{n} \\ &\equiv e^{-(\bar{w}.\bar{x} + \bar{b})} \bowtie \frac{1 - \bar{n}}{\bar{n}} \equiv \ln(e^{-(\bar{w}.\bar{x} + \bar{b})}) \bowtie \ln(\frac{1 - \bar{n}}{\bar{n}}) \\ &\equiv -(\bar{w}.\bar{x} + \bar{b}) \bowtie \ln(\frac{\bar{n}}{1 - \bar{n}}) \equiv \bar{w}.\bar{x} + \bar{b} \bowtie \ln(\frac{\bar{n}}{1 - \bar{n}}) \\ &\equiv \bar{x} \bowtie (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T).\ln(\frac{\bar{n}}{1 - \bar{n}}) - \bar{b}; [\text{From least square solution}] \\ &\equiv \bar{x} \bowtie (\bar{y}.\ln(\frac{\bar{n}}{1 - \bar{n}}) - \bar{b}); [\bar{y} = (\bar{w}^T.\bar{w})^{-1} . (\bar{w}^T)] \end{aligned}$$

1.4 Derivation of (BETATANH) rule:

For dense-based neural network with \tanh activation function, for given $N : y = \tanh(\overline{w}.\overline{x} + \overline{b})$ and $\delta : y \bowtie \overline{n}$, the computed precondition δ' is $\overline{x} \bowtie ((\overline{w}^T.\overline{w})^{-1}.\frac{1}{2}\ln(\frac{\overline{n}-1}{\overline{n}+1})) - \overline{b}$.

Given, postcondition $\delta : y \bowtie \overline{n}$ and given $N : y = \tanh(\overline{w}.\overline{x} + \overline{b})$ We know from the weakest precondition equation $\delta' = wp(N, \delta)$.

$$\begin{aligned}
 \text{So, } \delta' &= wp(N, \delta) \equiv wp(\tanh(\overline{w}.\overline{x} + \overline{b}), y \bowtie \overline{n}) \\
 &\equiv \tanh(\overline{w}.\overline{x} + \overline{b}) \bowtie \overline{n} \equiv \frac{e^{(\overline{w}.\overline{x} + \overline{b})} - e^{-(\overline{w}.\overline{x} + \overline{b})}}{e^{(\overline{w}.\overline{x} + \overline{b})} + e^{-(\overline{w}.\overline{x} + \overline{b})}} \bowtie \overline{n} \\
 \text{Let, } z &= \overline{w}.\overline{x} + \overline{b}. \text{ So, } P \equiv \frac{e^z - e^{-z}}{e^z + e^{-z}} \bowtie \overline{n} \\
 &\equiv \frac{2e^z}{-2e^{-z}} \bowtie \frac{\overline{n} + 1}{\overline{n} - 1} \equiv -\ln\left(\frac{e^z}{e^{-z}}\right) \bowtie \ln\left(\frac{\overline{n} + 1}{\overline{n} - 1}\right) \\
 &\equiv -(\ln(e^z) - \ln(e^{-z})) \bowtie \ln\left(\frac{\overline{n} + 1}{\overline{n} - 1}\right) \equiv -2z \bowtie \ln\left(\frac{\overline{n} + 1}{\overline{n} - 1}\right) \\
 &\equiv \overline{w}.\overline{x} + \overline{b} \bowtie -\frac{1}{2}.\ln\left(\frac{\overline{n} + 1}{\overline{n} - 1}\right) \equiv \overline{w}.\overline{x} + \overline{b} \bowtie \frac{1}{2}.\ln\left(\frac{\overline{n} - 1}{\overline{n} + 1}\right) \\
 &\equiv \overline{x} \bowtie ((\overline{w}^T.\overline{w})^{-1}.\frac{1}{2}\ln\left(\frac{\overline{n} - 1}{\overline{n} + 1}\right)) - \overline{b} \\
 &\equiv \overline{x} \bowtie (\overline{y}.\frac{1}{2}\ln\left(\frac{\overline{n} - 1}{\overline{n} + 1}\right) - \overline{b}); [\overline{y} = (\overline{w}^T.\overline{w})^{-1}.\overline{w}^T]
 \end{aligned}$$

1.5 Canonical Examples of Inferring Data Precondition from Derived wp rules

To infer data preconditions, starting from the last layer statement, we assert using the wp equation to determine the weakest precondition. The equation of *Dense* layer is as follows:

$$output = activation(dot(input, weight) + bias)$$

So, for given postcondition $\delta : y \leq n$, the statement S_3 can be written for output of last layer (y_3) with corresponding weight (\overline{w}_3) and bias(\overline{b}_3) as,

$$y_3 = \text{sigmoid}(\overline{w}_3^T.\overline{x}_3 + \overline{b}_3)$$

Let us look at a canonical example of how we can apply the rules presented in Fig. 4.

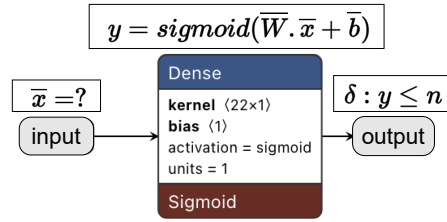


Figure 1: Example 1: Data precondition computation from given N with 1 layer and postcondition (δ)

1.5.1 Layer 1 Example: In Example 1 (Fig. 1), for given neural network (N) and postcondition (δ) as below,

$$N : a(f(\overline{x})) = \text{sigmoid}(\overline{W}.\overline{x} + \overline{b}); \delta : y \leq n$$

We can use wp rules over N and δ using (WPALPHA) rule to get the precondition for this single-layer neural network as follows,

$$wp(N, \delta) = wp(\text{sigmoid}(\overline{W}.\overline{x} + \overline{b}), \delta) \equiv \alpha(\delta, \beta(\text{sigmoid}(\overline{W}.\overline{x} + \overline{b})))$$

Now, we can apply (WPALPHASIGMA) and (BETASIGMOID) rules consecutively to get the data precondition as follows,

$$\begin{aligned}
 \alpha(\delta, \beta(\text{sigmoid}(\overline{W}.\overline{x} + \overline{b}))) &\equiv \beta(\text{sigmoid}(\overline{W}.\overline{x} + \overline{b}), \delta) \\
 &\equiv \beta(\text{sigmoid}(\overline{W}.\overline{x} + \overline{b}), y \leq n)
 \end{aligned}$$

$$\equiv \overline{x} \leq (\overline{y}.\ln\left(\frac{n}{1-n}\right)) - \overline{b}, \overline{y} = (\overline{W}^T.\overline{W})^{-1}.\overline{W}^T$$

Here, \overline{x} is an array of data that has been fed into the Neural network, and the predicate of \overline{x} denotes the precondition of the data for each feature that has been assumed by this 1-layer neural network.

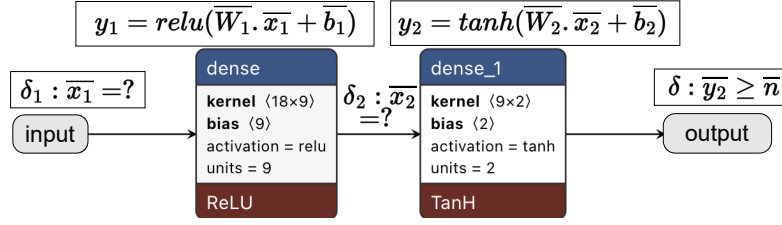


Figure 2: Example 2: Data precondition (δ_1) computation from given N with 2 layers and postcondition (δ)

1.5.2 Layer 2 Example: Let us consider, another model with different activation functions with 2 layers and multiple output classes (Example 2 in Fig. 2). Therefore, the output of this model (\bar{y}) is an array of n . Now, for this Example 2 Fig. 2, for given neural network (N) and postcondition (δ) as below,

$$N : a(f(\bar{x})).a(f(\bar{x})) = \text{relu}(\bar{W}_1.\bar{x}_1 + \bar{b}_1).\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2); \delta : \bar{y}_2 \geq \bar{n}$$

We can use wp rules over N and δ using (WP), (WPALPHA) rule to get the precondition for this multiple layer neural network as follows,

$$\begin{aligned} \delta_1 &= wp(N, \delta) = wp(N_0.N_1, \delta) \equiv \\ &wp(\text{relu}(\bar{W}_1.\bar{x}_1 + \bar{b}_1).\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \delta) \\ \delta_1 &= wp(\text{relu}(\bar{W}_1.\bar{x}_1 + \bar{b}_1), \delta_2); \delta_2 = wp(\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \delta) \end{aligned}$$

Now, we can apply (WPALPHASIGMA) and (BETATANH) rules consecutively to get the data precondition as follows,

$$\begin{aligned} \delta_2 &= wp(\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \delta) \equiv \alpha(\delta, \beta(\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2))) \\ &\equiv \beta(\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \delta) \equiv \beta(\text{tanh}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \bar{y}_2 \geq \bar{n}) \\ &\equiv \bar{x}_2 \leq (\bar{y}_2 \cdot \frac{1}{2} \ln(\frac{\bar{n}-1}{\bar{n}+1})) - \bar{b}_2, \bar{y}_2 = (\bar{W}_2^T.\bar{W}_2)^{-1} \cdot (\bar{W}_2^T) \end{aligned}$$

Here, \bar{x}_2 is an array of input that has been obtained from the first layer and fed into the second layer, and the predicate of \bar{x}_2 denotes the precondition of the data in layer 2, which is a postcondition of layer 1. After asserting with this postcondition, we obtain δ_1 similarly using (BETARELU) rule,

$$\begin{aligned} \delta_1 &= wp(\text{relu}(\bar{W}_1.\bar{x}_1 + \bar{b}_1), \bar{x}_2 \leq (\bar{y}_2 \cdot \frac{1}{2} \ln(\frac{\bar{n}-1}{\bar{n}+1})) - \bar{b}_2) \\ &\equiv \bar{x}_1 \leq (\bar{y}_1 \cdot (\bar{y}_2 \cdot \frac{1}{2} \ln(\frac{\bar{n}-1}{\bar{n}+1})) - \bar{b}_2) - \bar{b}_1 \wedge \bar{x}_1 \leq (\bar{y}_1 \cdot (-\bar{b}_1)), \\ &\bar{y}_1 = (\bar{W}_1^T.\bar{W}_1)^{-1} \cdot (\bar{W}_1^T) \end{aligned}$$

In this step, we obtain the precondition which is an array of the data (\bar{x}_1) for each feature that has been assumed by this neural network with 2 layers.

1.5.3 Layer 3 Example: Let us consider an example with multiple Dense layers with linear, sigmoid activation functions. Now, for this Example 3 Fig. 3, for given neural network (N) and postcondition (δ) as below, We have Dense layers with linear, linear, sigmoid activation functions for this example. Now, for given neural network (N) and postcondition (δ),

$$N : \text{linear}(\bar{W}_1.\bar{x}_1 + \bar{b}_1).\text{linear}(\bar{W}_2.\bar{x}_2 + \bar{b}_2).\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3);$$

$$\delta : y_3 \geq n_1 \wedge y_3 \leq n_2$$

Our proposed technique is generalized to DNN models with multiple layers. For example, a DNN model presented in Fig. 3 has 3 layers and different activation functions. In that model, the output layer has a single class, i.e., the output value $y \in \mathbf{R}$. The given postcondition is an instance of ($\delta \wedge \delta$) and will be in the range between $[n_1, n_2]$. Now, We utilize wp rules over N and δ using (WP), (WPALPHA) rules to get the precondition for this multiple-layer neural network as follows,

$$\begin{aligned} \delta_1 &= wp(N, \delta) = wp(N_0.N_1, \delta) \equiv \\ &wp(\text{linear}(\bar{W}_1.\bar{x}_1 + \bar{b}_1).\text{linear}(\bar{W}_2.\bar{x}_2 + \bar{b}_2).\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3), \delta) \\ \delta_1 &= wp(\text{linear}(\bar{W}_1.\bar{x}_1 + \bar{b}_1), \delta_2); \\ \delta_2 &= wp(\text{linear}(\bar{W}_2.\bar{x}_2 + \bar{b}_2).\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3), \delta); \\ \delta_2 &= wp(\text{linear}(\bar{W}_2.\bar{x}_2 + \bar{b}_2), \delta_3); \delta_3 = wp(\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3), \delta) \end{aligned}$$

Then, we apply (WPALPHASIGMA), (WPALPHAWEDGE), (BETASIGMOID) rules consecutively to get the precondition as follows,

$$\delta_3 = wp(\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3), \delta) \equiv \alpha(\delta_3, \beta(\text{sigmoid}(\bar{W}_3.\bar{x}_3 + \bar{b}_3)))$$

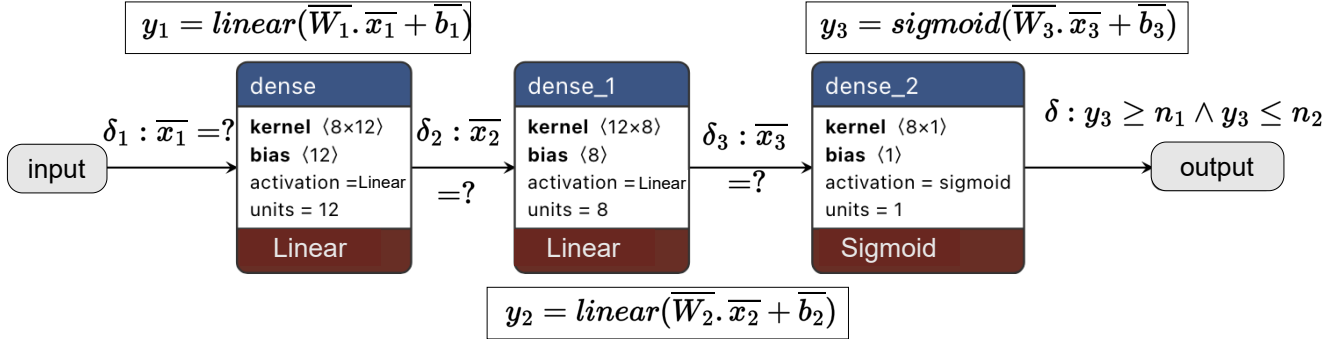


Figure 3: Data precondition (δ_1) computation from an example DNN model (N) with 3 layers and postcondition (δ)

$$\begin{aligned} &\equiv \beta(\text{sigmoid}(\overline{W}_3 \cdot \overline{x}_3 + \overline{b}_3), \delta_3) \equiv \beta(\text{sigmoid}(\overline{W}_3 \cdot \overline{x}_3 + \overline{b}_3), y_3 \geq n_1 \wedge y_3 \leq n_2) \\ &\equiv \overline{x}_3 \geq ((\overline{y}_3 \cdot \ln \frac{n_1}{1 - n_1}) - \overline{b}_3) \wedge \overline{x}_3 \leq ((\overline{y}_3 \cdot \ln \frac{n_2}{1 - n_2}) - \overline{b}_3) \end{aligned}$$

Here, \overline{x}_3 is an array of input that has been obtained from the second layer and fed into the third layer, and the predicate of \overline{x}_3 denotes the precondition of the data in layer 3, which is a postcondition of layer 2. Here, \overline{y}_3 is an inverse function of the layer's weight matrix (\overline{W}_3). Then, we obtain δ_2 similarly using the *wp* rules (WPALPHA), (WPALPHAWEDGE), (BETALINEAR) consecutively,

$$\begin{aligned} \delta_2 &= \text{wp}(\text{linear}(\overline{W}_2 \cdot \overline{x}_2 + \overline{b}_2), \delta_3) \\ &\equiv \alpha(\delta_3, \beta(\text{linear}(\overline{W}_2 \cdot \overline{x}_2 + \overline{b}_2))) \equiv \beta(\text{linear}(\overline{W}_2 \cdot \overline{x}_2 + \overline{b}_2), \delta_3) \\ &\equiv \beta(\text{linear}(\overline{W}_2 \cdot \overline{x}_2 + \overline{b}_2), \overline{x}_3 \geq ((\overline{y}_3 \cdot \ln \frac{n_1}{1 - n_1}) - \overline{b}_3) \wedge \\ &\quad \overline{x}_3 \leq ((\overline{y}_3 \cdot \ln \frac{n_2}{1 - n_2}) - \overline{b}_3)) \equiv \overline{x}_2 \geq ((\overline{y}_2 \cdot ((\overline{y}_3 \cdot \ln \frac{n_1}{1 - n_1}) - \overline{b}_3)) - \overline{b}_2) \wedge \\ &\quad \overline{x}_2 \leq ((\overline{y}_2 \cdot ((\overline{y}_3 \cdot \ln \frac{n_2}{1 - n_2}) - \overline{b}_3)) - \overline{b}_2), \overline{y}_2 = (\overline{W}_2^T \cdot \overline{W}_2)^{-1} \cdot (\overline{W}_2^T) \end{aligned}$$

In this step, we obtain the precondition which is an array of the input (\overline{x}_2) that has been obtained from the first layer and fed into the second layer, and the predicate of \overline{x}_2 denotes the precondition of the input in layer 2, which is a postcondition of layer 1. After asserting with this postcondition, we obtain δ_1 similarly using the *wp* rules (WPALPHA), (WPALPHAWEDGE), (BETALINEAR),

$$\begin{aligned} \delta_1 &= \text{wp}(\text{linear}(\overline{W}_1 \cdot \overline{x}_1 + \overline{b}_1), \delta_2) \equiv \alpha(\delta_2, \beta(\text{linear}(\overline{W}_1 \cdot \overline{x}_1 + \overline{b}_1))) \\ &\equiv \beta(\text{linear}(\overline{W}_1 \cdot \overline{x}_1 + \overline{b}_1), \delta_2) \\ &\equiv \overline{x}_1 \geq ((\overline{y}_1 \cdot ((\overline{y}_2 \cdot ((\overline{y}_3 \cdot \ln \frac{n_1}{1 - n_1}) - \overline{b}_3)) - \overline{b}_2) - \overline{b}_1) \wedge \\ &\quad \overline{x}_1 \leq ((\overline{y}_1 \cdot ((\overline{y}_2 \cdot ((\overline{y}_3 \cdot \ln \frac{n_2}{1 - n_2}) - \overline{b}_3)) - \overline{b}_2) - \overline{b}_1) \end{aligned}$$

Finally, we obtain the precondition, which is an array of the data (\overline{x}_1) for each feature that this DNN has assumed with multiple layers where $\overline{y}_1 = (\overline{W}_1^T \cdot \overline{W}_1)^{-1} \cdot (\overline{W}_1^T)$. In our proposed technique, the entire process of data precondition inference from a DNN model is automated and generalized for other models, which is performed during the training stage.