# Expressive power of Neural Networks

A boolean ~~function~~ involving $n$ variables maps from a domain $\{\pm 1\}^n$ to $\{\pm 1\}$ can be implemented by a neural network $H_{V,E,\sigma}$ architecture.

$(+1, -1)$

**Claim:** For every $n$, there exists a graph $(V, E)$ of depth 2 such that $H_{V,E,\text{sign}}$ contains all functions from $\{\pm 1\}^n$ to $\{\pm 1\}$.

2 layers # nodes $|V|$) ??

$\sigma : \text{sign}(\text{threshold})$

Domain

**Proof:** We construct a graph with
$$|V_0| = n+1$$
$$|V_1| = 2^n + 1$$
$$|V_2| = 1$$

$V_0$ : $n+1$   $V_1$   $V_2$ : O

Let $f : \{\pm 1\}^n \longrightarrow \{\pm 1\}$ be some **boolean function**.

We can **adjust the weights** so that the network will implement $f$.

Let $\underline{u_1} \cdots \underline{u_k}$ be all vectors in $\{\pm 1\}^n$ on which $f$ outputs 1.

Consider a vector $\underline{x}^{\text{bool}} \in \{\pm 1\}^n$

| | 1 2 ... n | f( ) |
|---|---|---|
| $\underline{u_1}$ | 1 1 · · · 1 | +1 ← 👈 |
| $2^n$ rows | 1 1 · · · -1 | -1 |
| $\underline{u_2}$ | -1 1 · · · -1 | +1 ← 👈 |

Note that if $\underline{x} \neq u_i$ then $\langle \underline{x}, \underline{u_i} \rangle \leq n-2$

if $\underline{x} = u_i$, then $\langle \underline{x}, \underline{u_i} \rangle = n$

Affine transform

$$\begin{array}{ccc} +1 & +1 & -1 \end{array} = 3$$
$$\begin{array}{ccc} +1 & +1 & -1 \end{array} = 1$$
$$\begin{array}{ccc} +1 & -1 & -1 \end{array}$$

$\therefore$ a function $g_i(\underline{x}) = \underset{\sigma}{\text{sign}}(\underbrace{\langle \underline{x}, u_i \rangle - n+1})$ equals $\boxed{\perp}$

bias.

if and only if $\underline{x} = \underline{u_i}$

$g_i(\underline{x})$ computational node   $\underline{x} \neq u_i$   $\text{sign}(\underset{\text{sign}(-1)=-1}{n-2-n+1})$   $n-n+1$   $\text{sign}(1)=1$

$V_1$

$u_1$
$x$ $u_2$
$u_3$

$g_1$ $\quad 1$

$g_2$ $\quad 1$ $\qquad$ OR operation $\begin{matrix}-1\\+1\\-1\end{matrix}$

$g_3$ $\quad 1$ $\quad V_2$

$O$

$\downarrow K$ $\quad O \; \underline{g_K}$

$O = 1$

$O = 1$

$O = 1$

$K$

$\begin{matrix} u_1 \\ u_2 \\ u_K \end{matrix}$ $\quad \begin{matrix} f \\ 1 \\ 1 \\ 1 \\ 1 \end{matrix}$

$2^n$

$\underline{u_1} \quad \underline{u_2} \quad \underline{u_3} \; \cdots \; \underline{u_K}$

$f(\quad) = 1$

XEnt
CE
loss

output $\underline{O}$

GT $y$

one-hot encoding.

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

$O_1$
$O_2$ prob vector
$O_3$

$y_1$
$y_2$ prob vector.
$y_3$

$\leftarrow$ Xent $\rightarrow$

K: #classes.

#output layer nodes = #classes

$O_K$

$y_K$

| | | |
|---|---|---|
| 100 | 0.2 | 3 |
| 101 | 0.21 | 94 $K=2$ |
| 98 | 0.19 | 6 |
| 97 | 0.19 | 5 |
| 74 | 0.16 | 3 |

confidence high.
is low

$\underline{\underline{\sigma}}$ (Affine).

vector of real numbers

$\begin{matrix} v_1 \\ \rightarrow v_2 \\ \vdots \\ v_K \end{matrix}$

softmax $\longrightarrow$

$\begin{matrix} O_1 \\ O_2 \\ \text{probability} \\ O_K \end{matrix}$

We can adapt the weights between $V_0$ and $V_1$ so that for every $i \in [k]$, the $i^{th}$ neuron is $g_i(\underline{x})$ i

The neuron in the output layer implements a disjunction i.e. OR of the functions $g_i(\underline{x})$.

AND
Bias $-k+1$

$$f(\underline{x}) = \text{sign}\left(\sum_{i=1}^{k} g_i(\underline{x}) + k - 1\right)$$

bias.

$-k+k-1 \geqslant -1$

$-(k-2)+k-1 \geqslant 1$

$\text{sign}(1)$

Even if we try to model functions of the form $\{0,1\}^n \to \{0,1\}$ the size of the network will be exponential in $n$.

A neural network can approximate 1-Lipschitz function $f: [-1,+1]^n \longrightarrow [-1,1]$ within a precision $\in$, but the size of the network will be exponential in $n$.

$f(u)$
$f(v)$

$n$-variables $[-1 \quad +1]$

$\|f(u) - f(v)\| \leqslant \|u-v\|^2$

$\longrightarrow [-1,+1]$

Softmax converts $k$ real valued predictions $v_1 \cdots v_k$ into output probabilities $o_1 \cdots o_k$ using the relation

$$O_i = \frac{\exp(v_i)}{\sum_{j=1}^{k} \exp(v_j)} \qquad \forall i \in \{1, \cdots k\}$$

The softmax is mostly paired with the cross-entropy loss. If the target probability distribution over the $k$-classes is given by the vector $y_1 \cdots y_k$ then the cross-entropy loss is defined as

$$L = -\sum_{i=1}^{k} y_i \log(o_i)$$