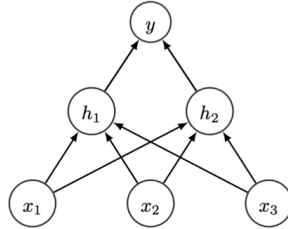# Mathematical Foundations of Machine Learning
## Assignment

**David Felipe Alvear Goyes**

### 1. Backpropagation.

The following graph shows the structure of a simple neural network with a single hidden layer. The input layer consists of three features $x = (x_1, x_2, x_3)$. The hidden layer includes two units $h = (h_1, h_2)$. The output layer includes one unit $\hat{y}$. We ignore bias terms for simplicity.



We use ReLU units $\varphi(z) = \max(0, z)$ as activation functions for the hidden and the output layer. Let $J = \ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$ be the loss function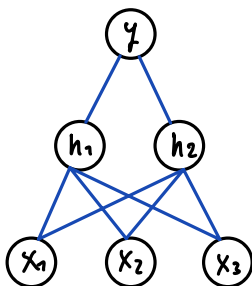. Here $y$ is the target value for the output unit. Denote by $W$ and $V$ weight matrices connecting input and hidden layer, and hidden layer and output respectively. They are initialized as follows:

$$W = \begin{bmatrix} 1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix} \text{ and } V = \begin{bmatrix} 0 & 1 \end{bmatrix}$$

- write out symbolically (in terms of $\varphi$, $W$, and $V$) the function $x \to \hat{y}$ represented by this network.
- consider a data sample with features $x = (1, 2, 1)$ and target $y = 1$. Compute the output value $\hat{y}$.
- compute the gradient of the loss function with respect to the weights. Specifically,

  - write down symbolically the gradient with respect to $V$, $\frac{\partial J}{\partial V}$
  - write down symbolically the gradient with respect to $W$, $\frac{\partial J}{\partial W}$
  - compute the values of this gradient numerically for the parameters above

## Solution

→ Write symbolically the function $x \to \hat{y}$



$h = \phi(w \cdot x) \longrightarrow h$ hidden layer

$\hat{y} = \phi(v \cdot h)$

$\hat{y} \to$ Predicted output

$\phi \to$ ReLU activation function

$\phi(z) = max(0, z)$

$\hat{y} = max\left(0, v\, max\left(0, w\, x\right)\right)$

$\hat{y} = \phi\left(v \cdot \phi(w \cdot x)\right)$

→ Considere $X = (1, 2, 1)$ $y = 1$, Compute $\hat{y}$.

In matrix form

$$\hat{y} = \text{Max}\left[0, \begin{bmatrix}0 & 1\end{bmatrix} \underbrace{\text{max}\left[0, \underbrace{\begin{bmatrix}1 & 0 & 1 \\ 1 & -1 & 0\end{bmatrix}}_{W} \underbrace{\begin{bmatrix}1 \\ 2 \\ 1\end{bmatrix}}_{X}\right]}_{\text{hidden layer}}\right]$$

where $V$ labels $\begin{bmatrix}0 & 1\end{bmatrix}$.

$$\hat{y} = \text{Max}\left[0, \underbrace{\begin{bmatrix}0 & 1\end{bmatrix}}_{V}\text{max}\left[0, \begin{bmatrix}2 \\ -1\end{bmatrix}\right]\right] = \text{Max}\left[0, \underbrace{\begin{bmatrix}0 & 1\end{bmatrix}}_{V}\begin{bmatrix}2 \\ 0\end{bmatrix}\right]$$

$$\hat{y} = \text{max}\left[0, \begin{bmatrix}0 \\ 0\end{bmatrix}\right] = 0 \quad\rightarrow\quad \boxed{\hat{y} = 0}$$

→ loss function: $J(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}$ $\qquad$ $\boxed{\dfrac{dJ}{d\hat{y}} = (\hat{y} - y) = -1}$

→ write $\dfrac{dJ}{dV} = \dfrac{dI}{d\hat{y}} \dfrac{d\hat{y}}{dV} = (\hat{y} - y)\,\phi'(V \cdot h)\,h$

→ write $\dfrac{dI}{dw} = \dfrac{dI}{d\hat{y}} \dfrac{d\hat{y}}{dh} \dfrac{dh}{dw} = (\hat{y} - y)\,\phi'(V \cdot h)\,V\,\phi'(w x)\,X$

$$h = \phi\left[\begin{bmatrix}1 & 0 & 1 \\ 1 & -1 & 0\end{bmatrix}\begin{bmatrix}1 \\ 2 \\ 1\end{bmatrix}\right] = \phi\begin{bmatrix}2 \\ -1\end{bmatrix} = \begin{bmatrix}2 \\ 0\end{bmatrix}$$

$$\hat{y} = \phi\left[\begin{bmatrix}0 & 1\end{bmatrix}\begin{bmatrix}2 \\ 0\end{bmatrix}\right] = 0$$

→ Numerically: $\dfrac{dI}{dV} = (\hat{y} - y)\,\phi'(V \cdot h)\,h = (-1)\,(0)\begin{bmatrix}2 \\ 0\end{bmatrix} = \begin{bmatrix}-2 \\ 0\end{bmatrix}$

$$\dfrac{dI}{dw} = (\hat{y} - y)\,\phi'(V \cdot h)\,V\,\phi'(w x)\,X = 0$$

**5. Sentiment analysis with trainable embeddings**

Instead of using fixed (pre-trained) word embeddings for the sentiment classification task as we did in the previous exercise, we will now make the word embeddings trainable as well. We will use an embedding of dimension $D = 50$.

- the classification network for this problem consists, as above, of an embedding layer, an averaging operation (with no trainable variables), a dense hidden layer of size 32, and a final dense output layer. Write down the mathematical description of the operations involved in evaluating the function defined by this network.

$$E = \frac{1}{N} \sum_{i=1}^{N} E_i \quad \longrightarrow \text{Embedding layer} + \text{averaging operation}$$

$$H = \emptyset \left( W_1 \cdot E + b_1 \right) \longrightarrow \text{Hidden layer} \qquad W_1^{32 \times 50}$$

$$\hat{y} = W_2 \cdot H + b_2 \longrightarrow \text{output} \qquad W_2^{1 \times 32}$$