# Feed forward networks

| $x_1$ | $x_2$ | $t$ |
|----|----|----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$4$ examples

$$\sum_{i=1}^{4} \left( (w_1 x_1^{(i)} + w_2 x_2^{(i)} + b) - t_i \right)^2$$

Linear regression $\ne$ neural network

$w_1 = 0 \quad w_2 = 0 \ne b = 0.5$

→ matured features



$$h = f^{(1)}(x; W, C)$$
$$y = f^{(2)}(h, w, b) = w^T h + b$$

$$y = f^{(2)} f^{(1)} (x; W, C, w, b)$$

What will be $f^{(1)}$?

L.T. $\begin{bmatrix} 2 & 1 \end{bmatrix}$

$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \in \mathbb{R}^2$ to $\mathbb{R}^3$ can be $\begin{bmatrix} f_1(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \\ f_2(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \\ f_3(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \end{bmatrix}$

$\begin{bmatrix} 3 & 1 & 5 \\ 4 & 2 & 6 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 10 \\ ? \\ 16 \end{bmatrix}$

$f_1 = 3x_1 + 4x_2$

$f_2 = x_1 + 2x_2$ } linear combination

$f_3 = 5x_1 + 6x_2$ → linear transformation

# Affine transformation

Affine transform: Linear transformation + constant vector added to result

$$f^{(1)}(x: W, C) = \underbrace{W^T x + b}_{\text{Affine transform}} \quad \} \text{matured space}$$

For brevity (and to make life easier), let us just consider a linear transformation

$$y = \omega^T (W^T x)$$
$$= (\omega')^T x \quad \text{where} \quad (\omega')^T = \omega^T W^T$$

$$\sum_{i=1}^{4} (\omega^T x - t_i)^2$$

∴ SO, The resultant matured features of a linear transformation is a linear model. (Which can't fit a non linear function like XOR)
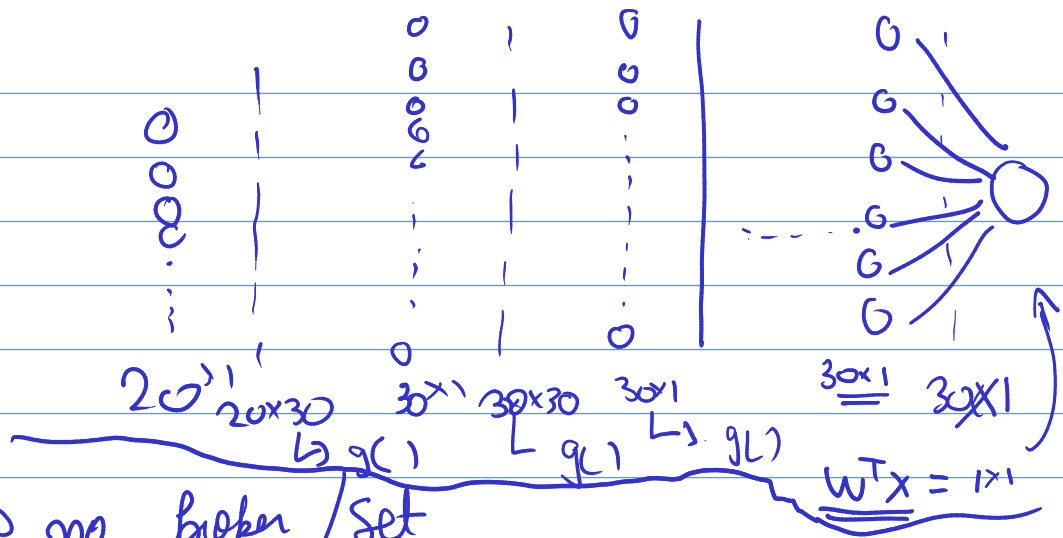
So as a result of this problem, all of deep learning and
Some machine learning algorithms use something else on top of it
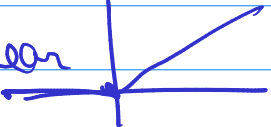
$$\begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} + \begin{bmatrix} C_1 \\ C_2 \end{bmatrix} = \begin{bmatrix} W_{11} + W_{21} + C_1 \\ W_{12} + W_{22} + C_2 \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}$$

Add a non linear function $= \begin{bmatrix} g(W_{11} + W_{21} + C_1) \\ g(W_{12} + W_{22} + C_2) \end{bmatrix}$

one of them could be

$g(z) = \max(z, 0) \Rightarrow$ This function is an activation function

applied to every hidden layer

$$20^{\times 1} \quad 20\times30 \quad 30^{\times 1} \quad 30\times30 \quad 30\times 1 \qquad 30\times 1 \quad 30\times 1$$

$$\underset{g()}{\hookrightarrow} \qquad \underset{g()}{\llcorner} \qquad \underset{g(L)}{\llcorner}$$

$$\underline{\underline{W^T x = 1\times 1}}$$

There's no proper / set "algorithm" to choose an activation function, You have to experiment / make a judgement to choose

RELU: piecewise linear  but __not__ actually linear → nonlinear

→ It gives nonlinear properties, and can be used as activation
  ⇒ but you need to do the math involved.
→ It even has some biological motivation.

Similarly, there is not set algo for number of layers
    we will needs insights for the domain/
    experimentation if to the ~~above insights~~ are not
    there.