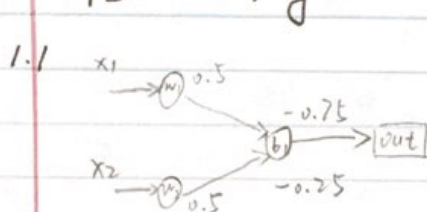


PS2 Yang Tian

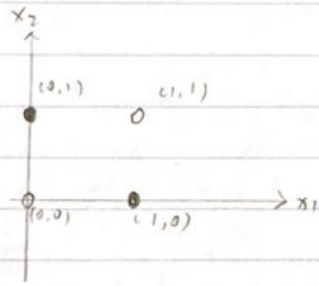


$$W_{AND} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad b_{AND} = -0.75$$

$$W_{OR} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad b_{OR} = -0.25$$

1.2

x_1	x_2	$f_{XOR}(x)$
0	0	0
0	1	1
1	0	1
1	1	0



$f(x) = W^T x + b$
 is a linear model.
 from left figure, we
 can find that no
 trial linear separator (line)
 can separate two classes.

2.1 Since $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable and Lipschitz continuous

$$|f(x_1) - f(x_0)| \leq k \|x_1 - x_0\|$$

$$\frac{\|f(x_1) - f(x_0)\|}{\|x_1 - x_0\|} \leq k$$

According to mean value theorem, for $\forall x, y \in [a, b]$, $x < y$ $\exists c \in (x, y) \subseteq (a, b)$
 s.t.

$$\left| \frac{f(y) - f(x)}{y - x} \right| = |f'(c)| \leq L$$

Furthermore, we could argue that Lipschitz constant $k \geq \sup_{x \in (x_0, x_1)} |f'(x)|$
 Therefore $|\nabla f(x_0)| \leq k$

$$|R(x)| = |f(x) - f(x_0) - \nabla f(x_0)^T (x - x_0)|$$

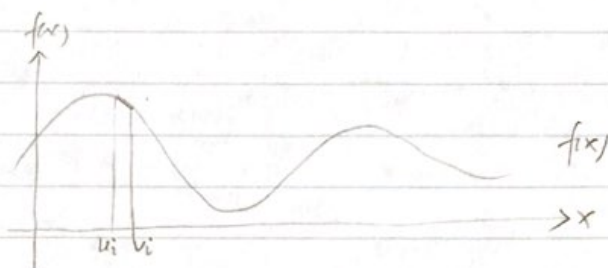
$$\leq |f(x) - f(x_0)| + |\nabla f(x_0)^T (x - x_0)|$$

$$\leq k \|x - x_0\| + k \|x - x_0\| \leq 2kC$$

2.2. Activation function

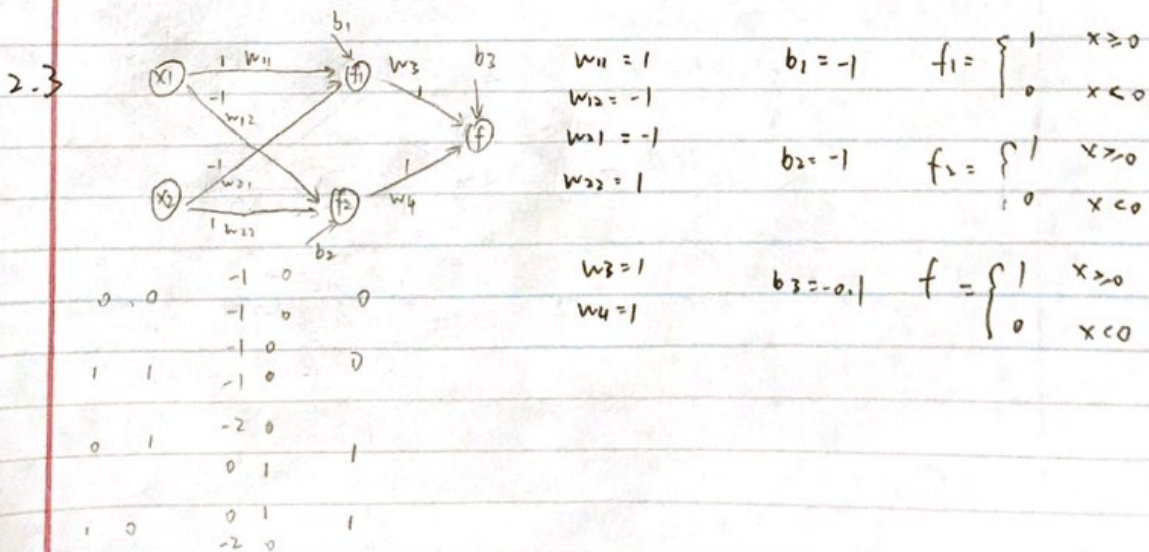
$$g(h, u, v)_i = \begin{cases} 0 & h_i < u_i \\ h_i & u_i < h_i < v_i \\ 0 & v_i < h_i \end{cases}$$

From the figure, we know it is a bump function, whose bumper slope could further adjusted by $W^{(2)}$. Using a bunch of these functions, we could approximate the differentiable and Lipschitz continuous function $f(x)$ to any desired error ϵ .



For better approximation ~~and~~, we could use a larger but finite number of hidden units. Each unit only approximates $f(x)$, $\forall x \leq v_i$. Therefore the number of hidden units and corresponding weights ~~depend~~ depend on the total range of x and approximation error.

Theoretically speaking we could think of approximation as wavelet transform. The activation function is mother wave. Using wavelet transform, we could find the parameters for each wave and the total number of ~~needed~~ waves needed, given the desired approximation error.



$$3. \quad h(x) = W^{(3)} \max \{ 0, W^{(2)} \max \{ 0, W^{(1)} x + b^{(1)} \} + b^{(2)} \} + b^{(3)}$$

$$3.1 \quad x=1$$

$$h(x) = W^{(3)} \max \{ 0, W^{(2)} \cdot \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} + b^{(2)} \} + b^{(3)}$$

$$= W^{(3)} \max \{ 0, \begin{bmatrix} 2 \\ 2 \end{bmatrix} \} + b^{(3)}$$

$$= 5$$

$$\frac{dh}{dx} = [1, 1] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = 2$$

$$h(x) = wx + b$$

$$w=2 \quad b=3$$

$$3.2 \quad x=-1$$

$$h(x) = W^{(3)} \max \{ 0, W^{(2)} \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} + b^{(2)} \} + b^{(3)}$$

$$= W^{(3)} \max \{ 0, \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \} + b^{(3)}$$

$$= 2$$

$$\frac{dh}{dx} = [1, 1] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = 1$$

$$h(x) = wx + b$$

$$w=1 \quad b=3$$

$$3.3 \quad x=-0.5$$

$$h(x) = W^{(3)} \max \{ 0, W^{(2)} \begin{bmatrix} 0 \\ 0.75 \end{bmatrix} + b^{(2)} \} + b^{(3)}$$

$$= W^{(3)} \max \{ 0, \begin{bmatrix} 0.75 \\ 0.75 \end{bmatrix} \} + b^{(3)}$$

$$= 2.5$$

$$\frac{dh}{dx} = 1$$

$$h(x) = wx + b$$

$$w=1 \quad b=3$$

$$4. (a) \quad W = 2 \cdot I = \begin{bmatrix} 2 & & 0 \\ & \ddots & \\ 0 & & 2 \end{bmatrix} \quad b = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

For d -dimension space, there are 2^d regions of input are identified onto $0 = (0, 1)^d$ by $f(x) = |W^{(1)}x + b^{(1)}|$

when $d=1$ we can easily find $R_1 = (0.5, 1)$ and $R_2 = (0, 0.5)$

when $d=2$ we need to choose either R_1 or R_2 for each dimension. so there are $2^2 = 4$ possibility regions

Generally speaking, for d -dimension, there are two possibilities for each dimension, so there are 2^d regions.

(b) Without loss of ~~gener~~ generality, we could assume $g(\cdot)$ is applied on input first, then $f(\cdot)$.

For function $f(\cdot)$, there are n_f regions that could be identified onto $(0, 1)^d$, each of which, obtained by applying $g(\cdot)$ on input field, has n_g identified ~~reg~~ regions.

As a result, totally speaking there are $n_f \cdot n_g$ regions identified onto final $[0, 1]^d$.

(c) Let's start backwards from layer L .

Output region is $[0, 1]^d$. So h_L will have 2^d identified regions. Normally not all of them (or even ~~none~~ none of them) are still $[0, 1]^d$. But the result from question (a) doesn't have any restriction on output region. So for each of 2^d regions before layer L , there are 2^d regions at layer $L-1$. And this goes on and on, we will find down to the beginning, there are 2^{Ld} regions identified to $[0, 1]^d$, which is exponential to the depth of networks and dimension of space.