

1.

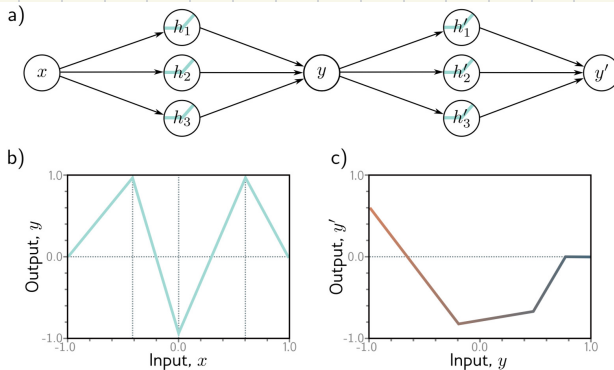
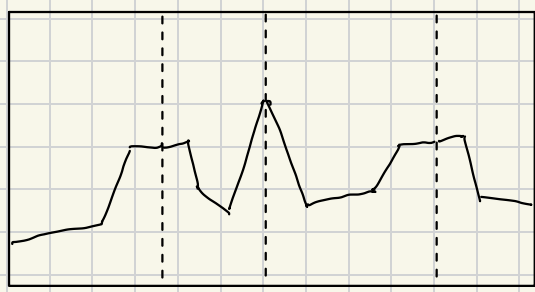


Figure 4.8 Composition of two networks for problem 4.1. a) The output y of the first network becomes the input to the second. b) The first network computes this function with output values $y \in [-1, 1]$. c) The second network computes this function on the input range $y \in [-1, 1]$.



3. Non-negative homogeneity of ReLU:

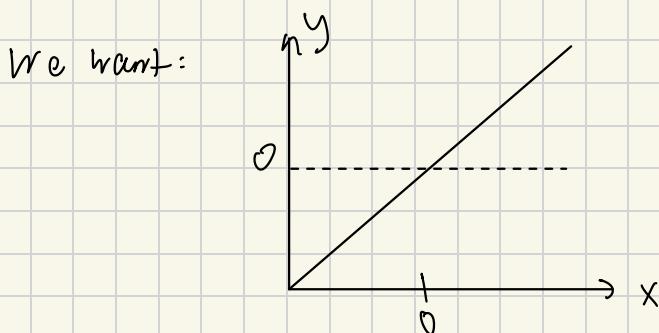
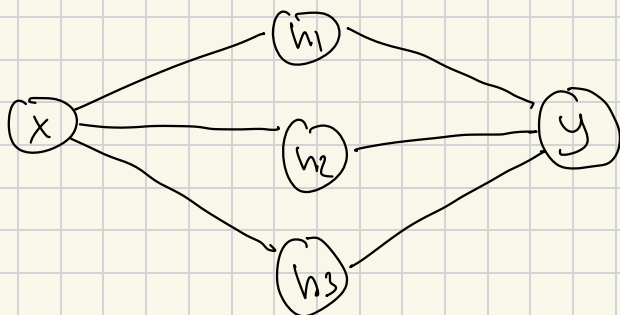
$\star \text{ReLU}(\alpha z) = \alpha \text{ReLU}(z)$

$$\begin{aligned}
 & \text{ReLU}[\beta_1 + \gamma_1 \cdot z, \text{ReLU}[\beta_0 + \gamma_0 \cdot z_0 x]] \\
 &= \text{ReLU}\left[\gamma_0 \gamma_1 \left(\frac{\beta_1}{\gamma_0 \gamma_1} + \frac{\gamma_1}{\gamma_0 \gamma_1} \text{ReLU}\left[\gamma_0 \left(\frac{\beta_0}{\gamma_0} + \frac{\gamma_0}{\gamma_0} z_0 x \right) \right] \right)\right] \\
 &= \gamma_0 \gamma_1 \text{ReLU}\left[\frac{\beta_1}{\gamma_0 \gamma_1} + \frac{\gamma_1}{\gamma_0} \gamma_0 \text{ReLU}\left[\frac{\beta_0}{\gamma_0} + z_0 x \right] \right] \\
 &= \gamma_0 \gamma_1 \text{ReLU}\left[\frac{\beta_1}{\gamma_0 \gamma_1} + \gamma_1 \text{ReLU}\left[\frac{\beta_0}{\gamma_0} + z_0 x \right] \right]
 \end{aligned}$$

5. Depth \rightarrow No layers = 20
 width \rightarrow No units / layer = 30

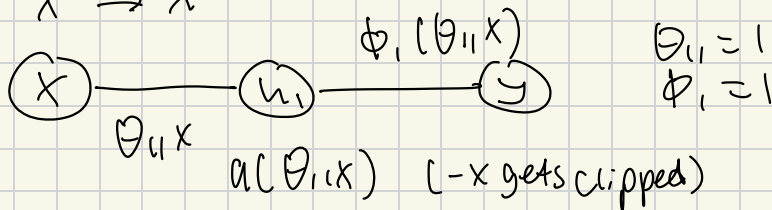
7.
$$y = f[x, \phi]$$

$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]. \quad (3.1)$$

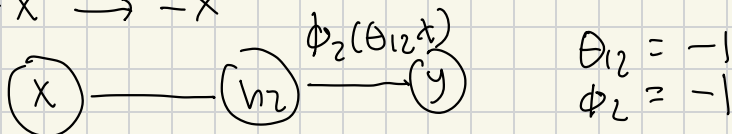


$$\begin{aligned} -x &\rightarrow -x \\ x &\rightarrow x \end{aligned}$$

(1) $x \rightarrow x$



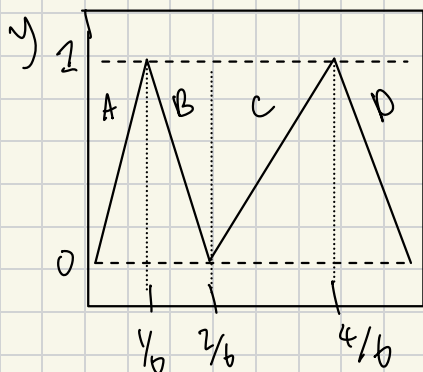
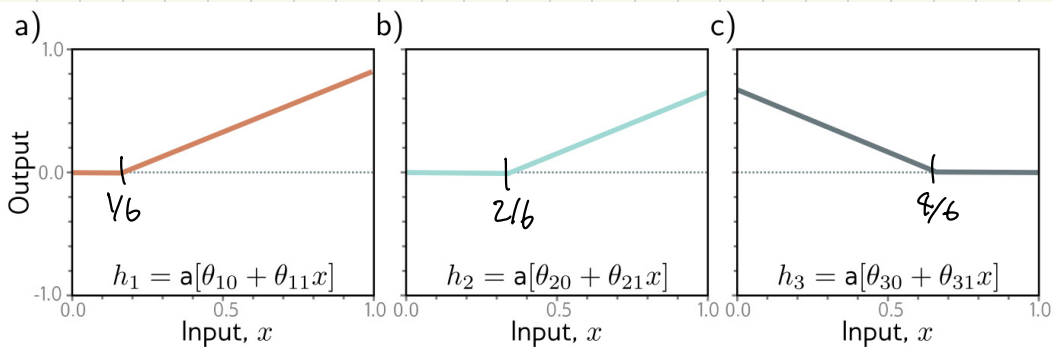
(2) $-x \rightarrow -x$



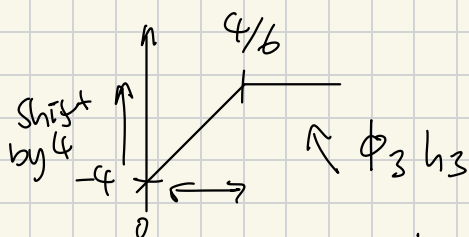
[other params
can be
set to zero]

$y = \text{ReLU}[\theta_{12}x] + \text{ReLU}[-x]$ (x gets clipped)

8.



$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$



A: only h_3 is active. $x=0, y=0$, $x=1/6, y=2 \rightarrow \phi_3 = -6$

B: h_1, h_3 active. Slope $\phi_1 - \phi_3$

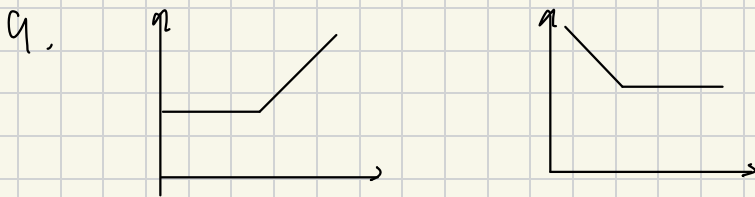
We need a slope of $-6 = \phi_1 + 6 \rightarrow \phi_1 = -12$

C: All active. Slope $\phi_1 + \phi_2 - \phi_3$

We need a slope of $3 = -12 + \phi_2 + 6$
 $3 = -6 + \phi_2$
 $\rightarrow \phi_2 = 9$

$\therefore \phi_0 = 4, \phi_1 = -12, \phi_2 = 9, \phi_3 = -6$

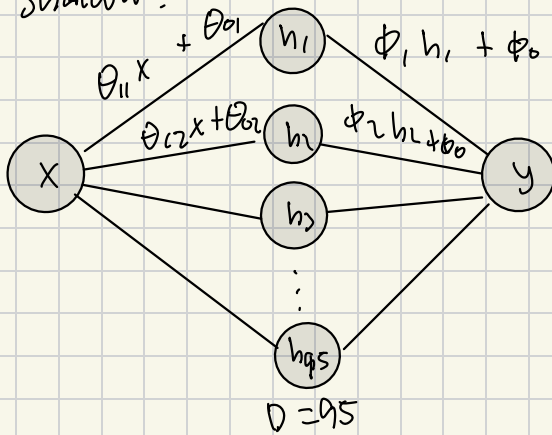
- Comparing the network on itself creates 4 new regions in each region $\rightarrow 4 \times 4 = 16$ new regions. $= 4^2$
- Composing this k times is 4^{k+1} new regions.



Imagine activations of two hidden units

- No matter how you scale and add them, you cannot create three linear regions that oscillate between 0 and 1.
- But this is possible to create 5 linear regions oscillating between 0 and 2 with 4 units (not shown).

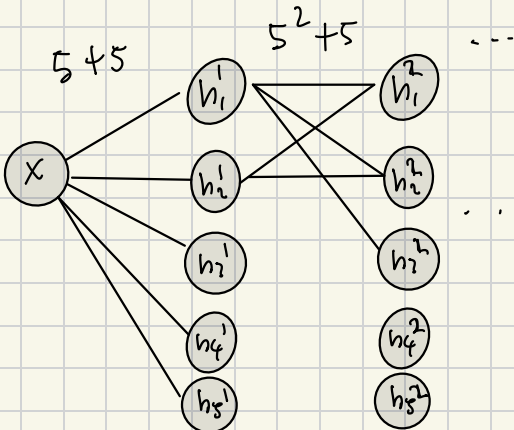
11. Shallow:



Contains

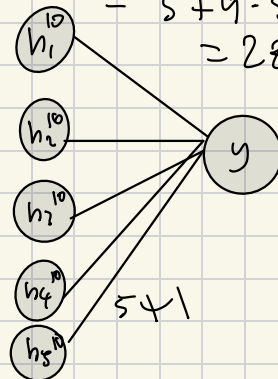
- 95×2 weights
 - $95 + 1$ biases
- $= 286$ parameters.

Deep:



Contains

- $5 + 9 \cdot 5^2 + 5$ weights
 - $5 + 9 \cdot 5 + 1$ biases
- $= 286$ parameters



Shallow network has 46 linear regions
- Deep network has

$$N_r = \left(\frac{D}{D_i} + 1\right)^{D_i(k-1)} \sum_{j=0}^{D_i} \binom{D}{j}$$

$$D_i = 1, D = 5, k = 10$$

$$N_r = (5 + 1)^9 = 6^9$$

Typically, the shallow network will run faster because the operations are parallelizable, whereas in the deep network, the outputs of each layer must be mapped into the following layers (sequential).