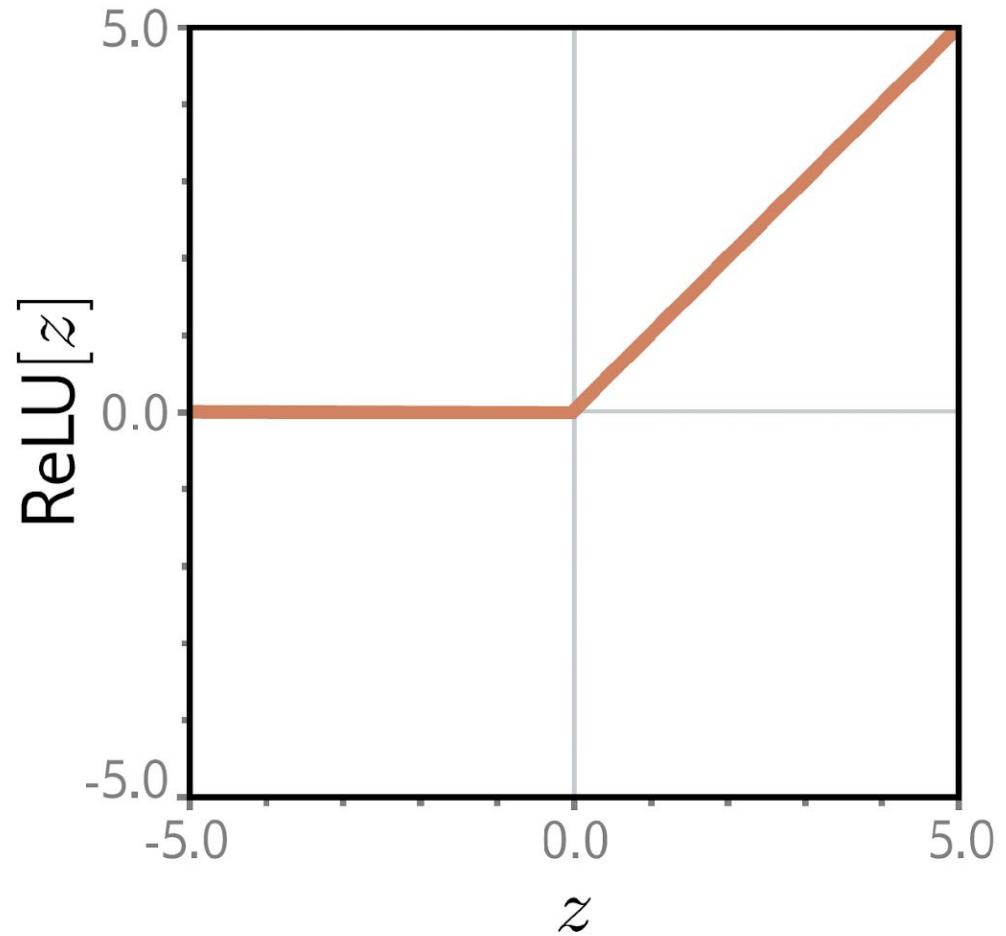
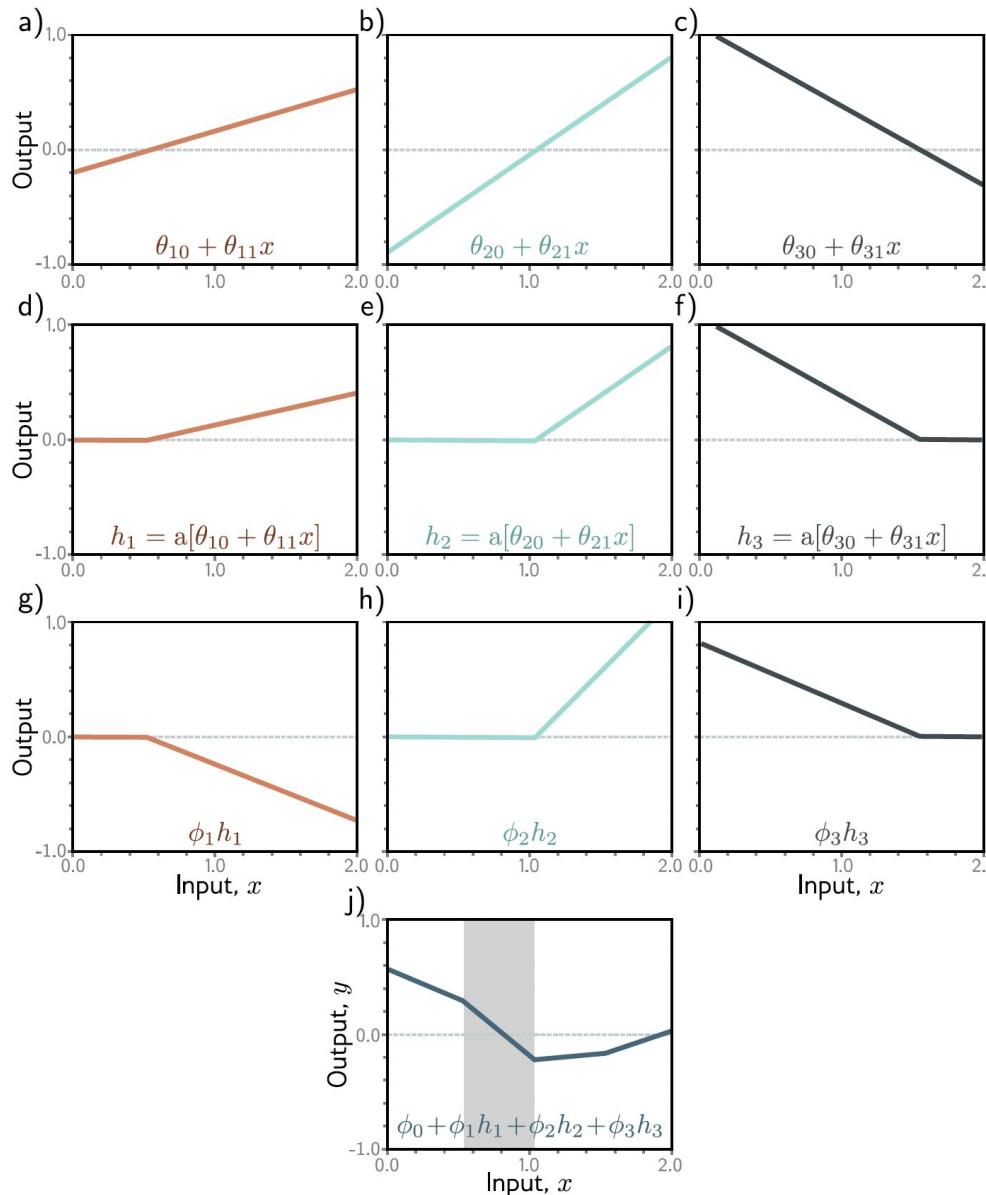


# Understanding Deep Learning

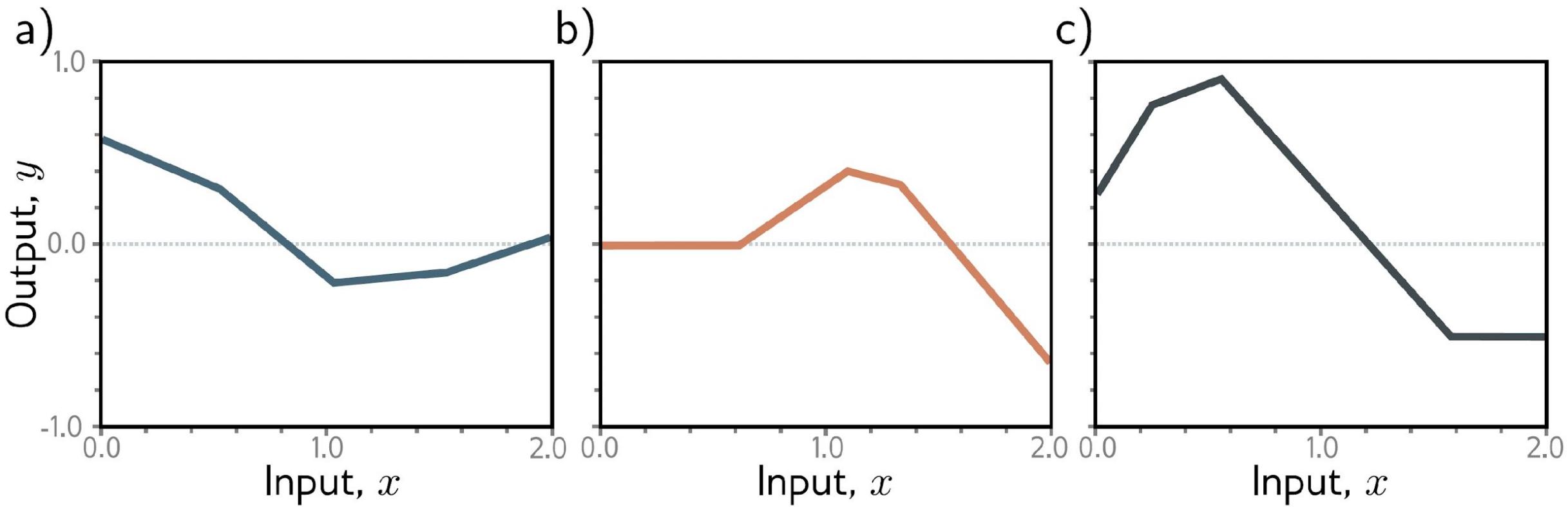
Chapter 3: Shallow Neural Networks

**Figure 3.1** Rectified linear unit (ReLU). This activation function returns zero if the input is less than zero and returns the input unchanged otherwise. In other words, it clips negative values to zero. Note that there are many other possible choices for the activation function (see figure 3.13), but the ReLU is the most commonly used and the easiest to understand.

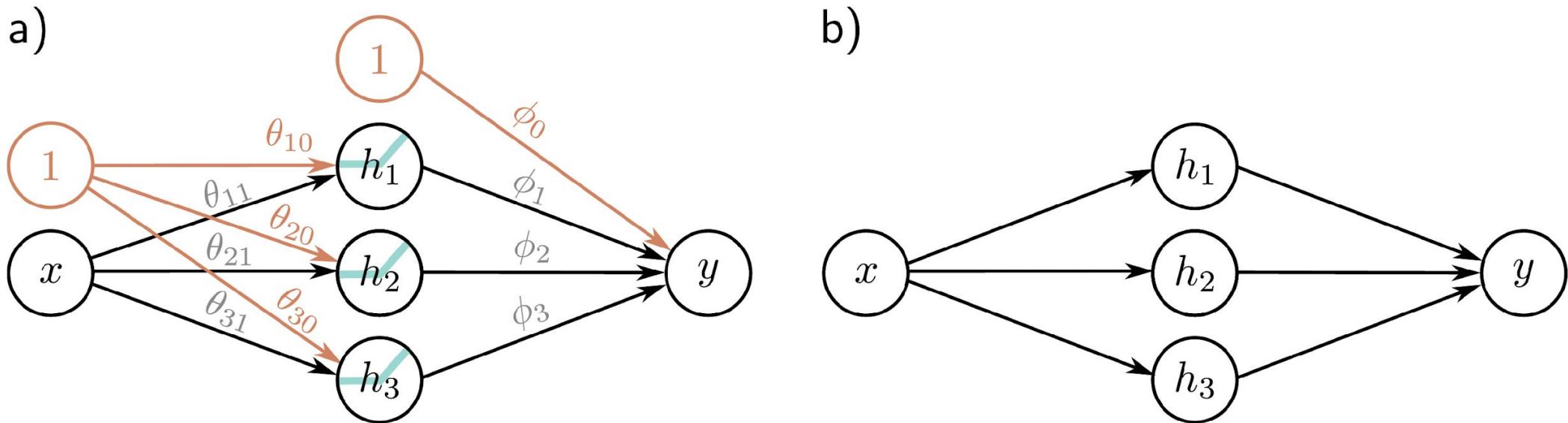




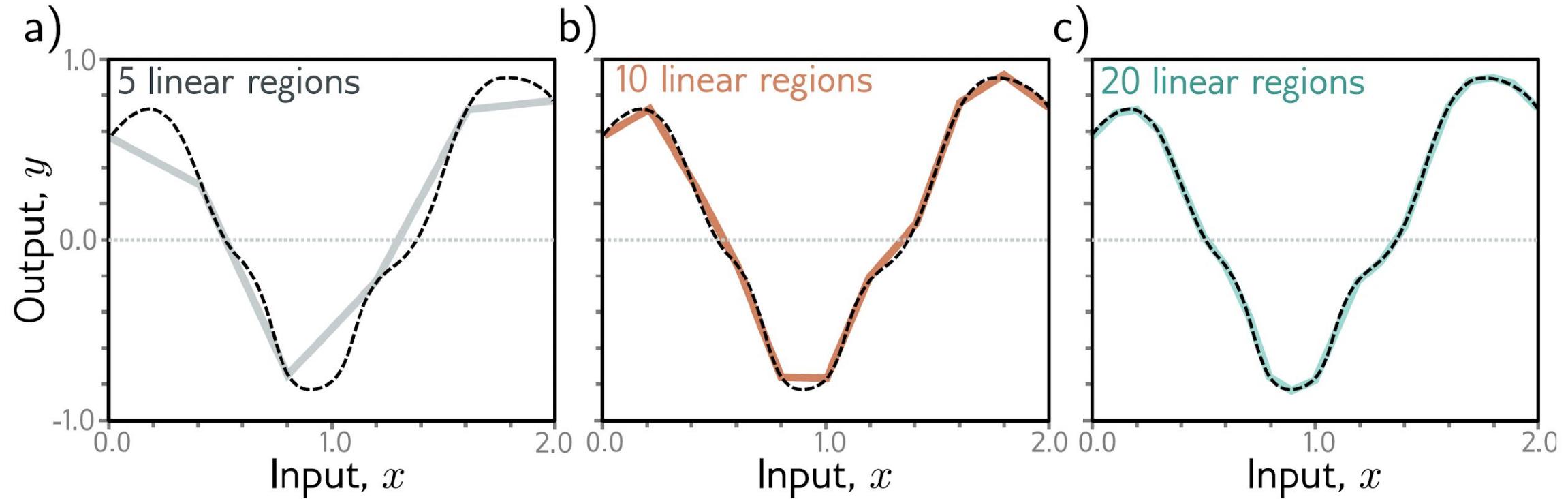
**Figure 3.3** Computation for function in figure 3.2a. a–c) The input  $x$  is passed through three linear functions, each with a different y-intercept  $\theta_{\bullet 0}$  and slope  $\theta_{\bullet 1}$ . d–f) Each line is passed through the ReLU activation function, which clips negative values to zero. g–i) The three clipped lines are then weighted (scaled) by  $\phi_1, \phi_2$ , and  $\phi_3$ , respectively. j) Finally, the clipped and weighted functions are summed, and an offset  $\phi_0$  that controls the height is added. Each of the four linear regions corresponds to a different activation pattern in the hidden units. In the shaded region,  $h_2$  is inactive (clipped), but  $h_1$  and  $h_3$  are both active.



**Figure 3.2** Family of functions defined by equation 3.1. a–c) Functions for three different choices of the ten parameters  $\phi$ . In each case, the input/output relation is piecewise linear. However, the positions of the joints, the slopes of the linear regions between them, and the overall height vary.

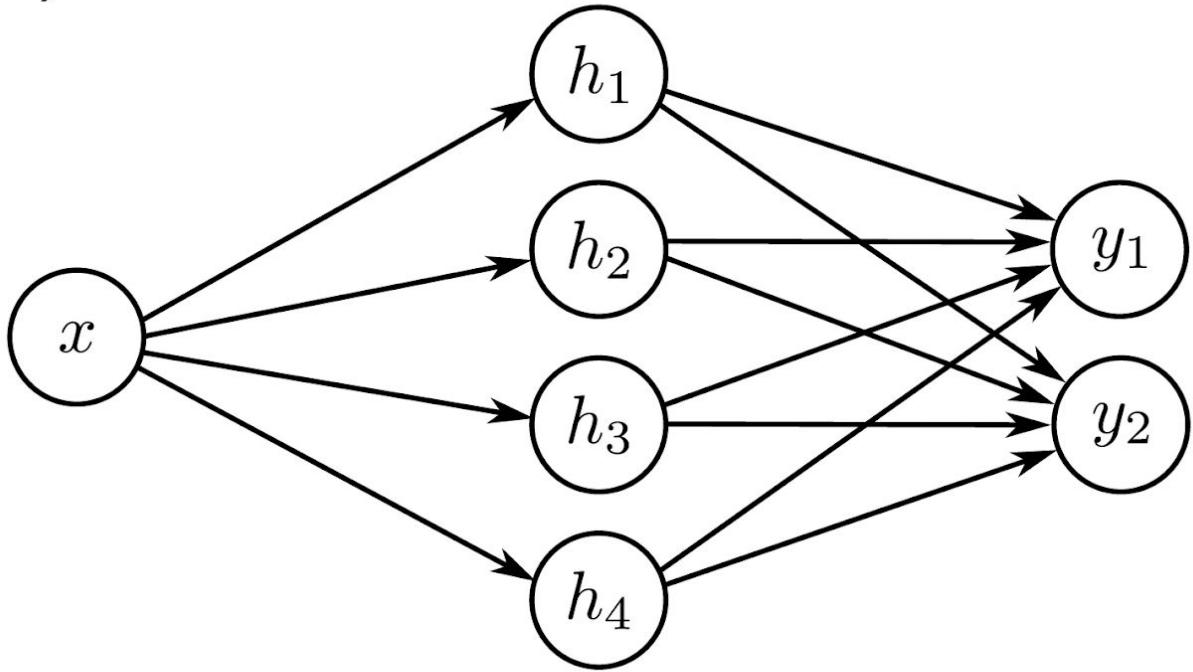


**Figure 3.4** Depicting neural networks. a) The input  $x$  is on the left, the hidden units  $h_1$ ,  $h_2$ , and  $h_3$  in the center, and the output  $y$  on the right. Computation flows from left to right. The input is used to compute the hidden units, which are combined to create the output. Each of the ten arrows represents a parameter (intercepts in orange and slopes in black). Each parameter multiplies its source and adds the result to its target. For example, we multiply the parameter  $\phi_1$  by source  $h_1$  and add it to  $y$ . We introduce additional nodes containing ones (orange circles) to incorporate the offsets into this scheme, so we multiply  $\phi_0$  by one (with no effect) and add it to  $y$ . ReLU functions are applied at the hidden units. b) More typically, the intercepts, ReLU functions, and parameter names are omitted; this simpler depiction represents the same network.

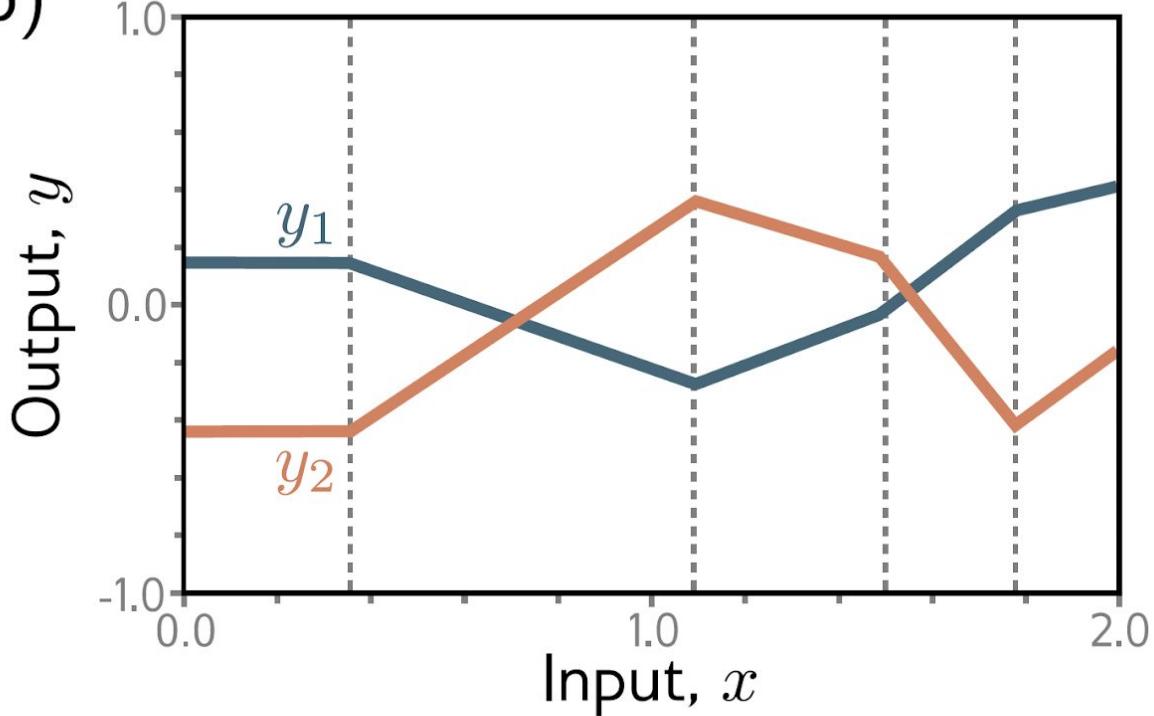


**Figure 3.5** Approximation of a 1D function (dashed line) by a piecewise linear model. a–c) As the number of regions increases, the model becomes closer and closer to the continuous function. A neural network with a scalar input creates one extra linear region per hidden unit. The universal approximation theorem proves that, with enough hidden units, there exists a shallow neural network can describe any given continuous function defined on a compact subset of  $\mathbb{R}^D$  to arbitrary precision.

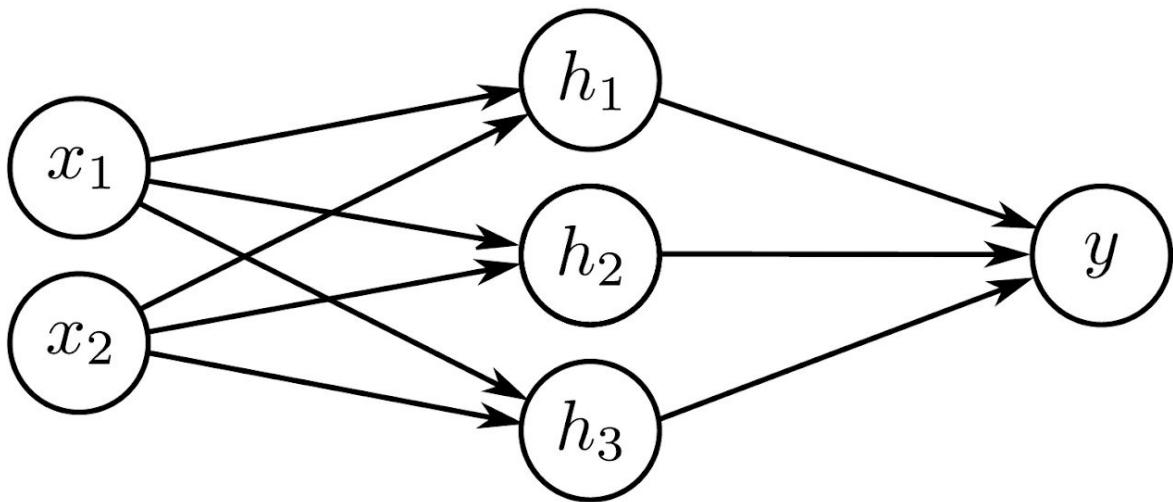
a)



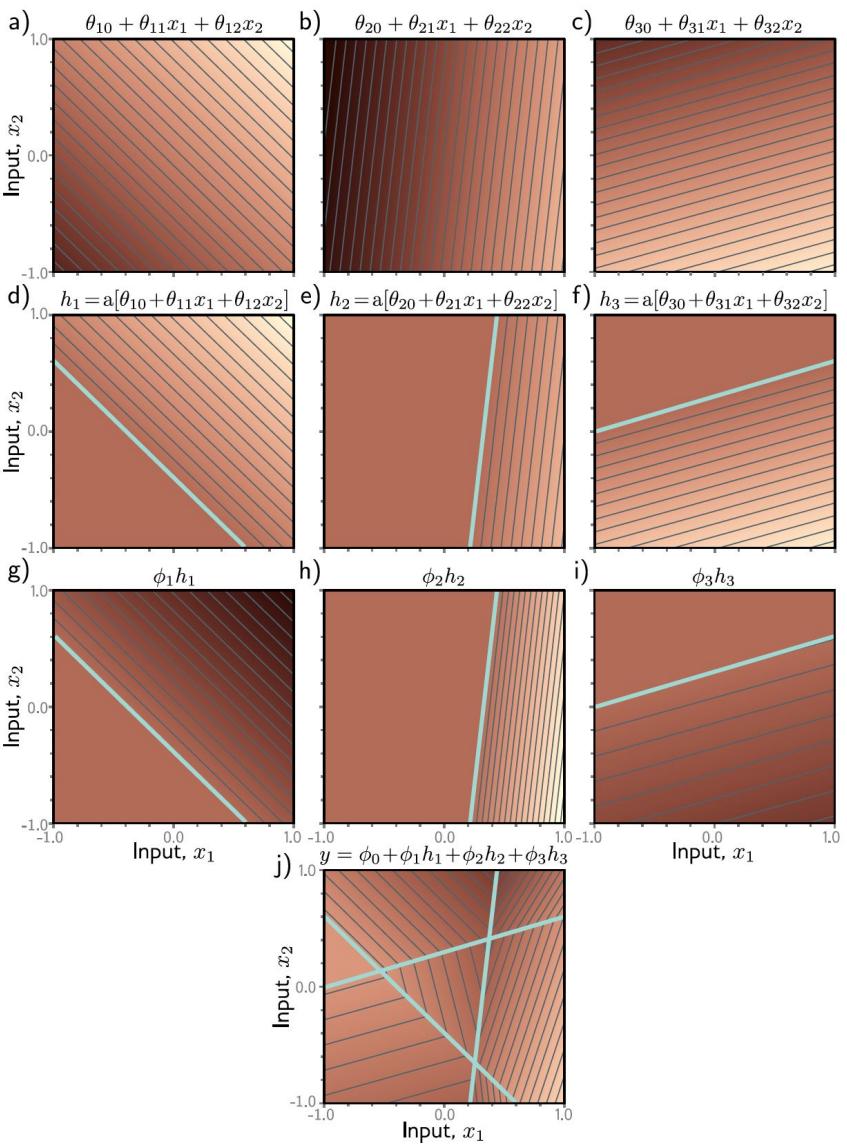
b)



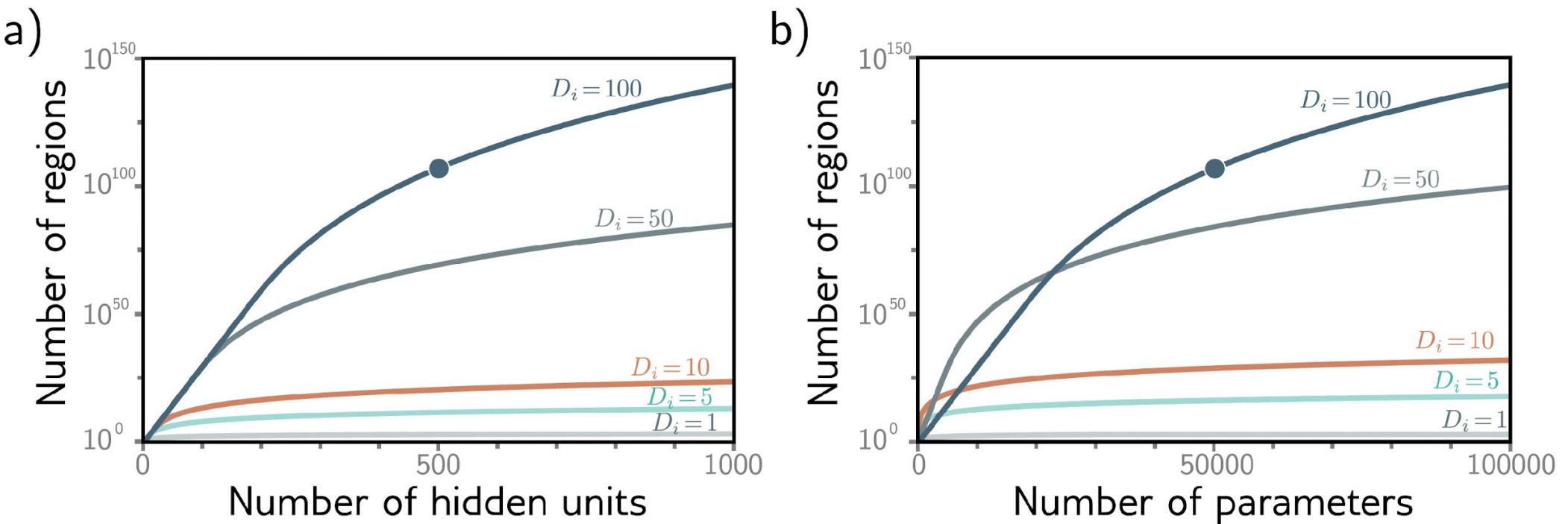
**Figure 3.6** Network with one input, four hidden units, and two outputs. a) Visualization of network structure. b) This network produces two piecewise linear functions,  $y_1[x]$  and  $y_2[x]$ . The four “joints” of these functions (at vertical dotted lines) are constrained to be in the same places since they share the same hidden units, but the slopes and overall height may differ.



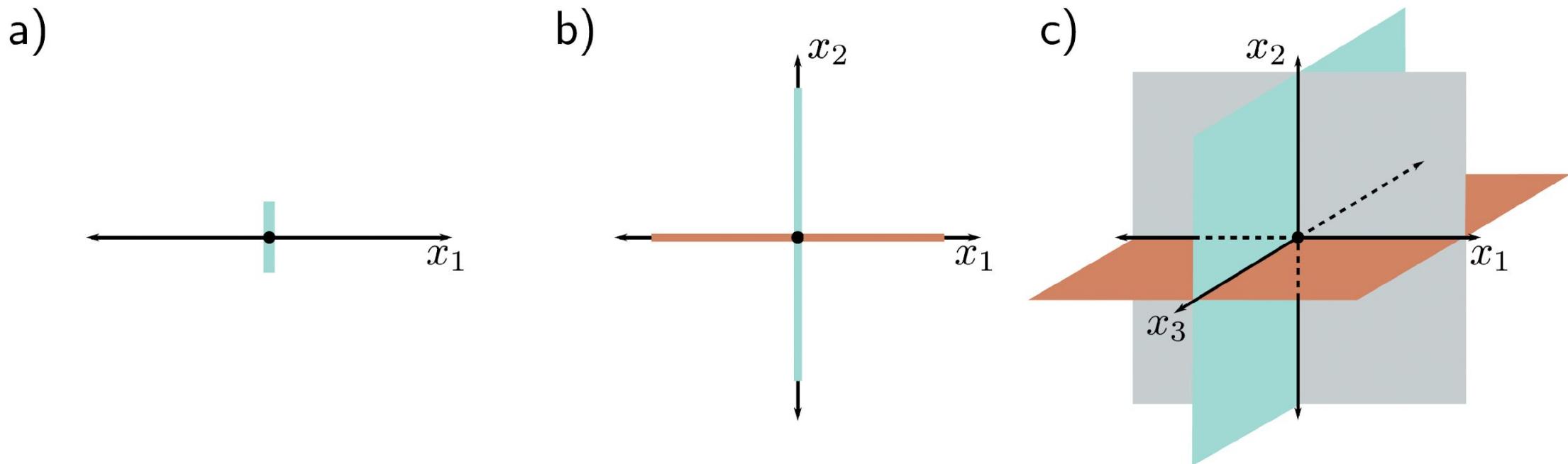
**Figure 3.7** Visualization of neural network with 2D multivariate input  $\mathbf{x} = [x_1, x_2]^T$  and scalar output  $y$ .



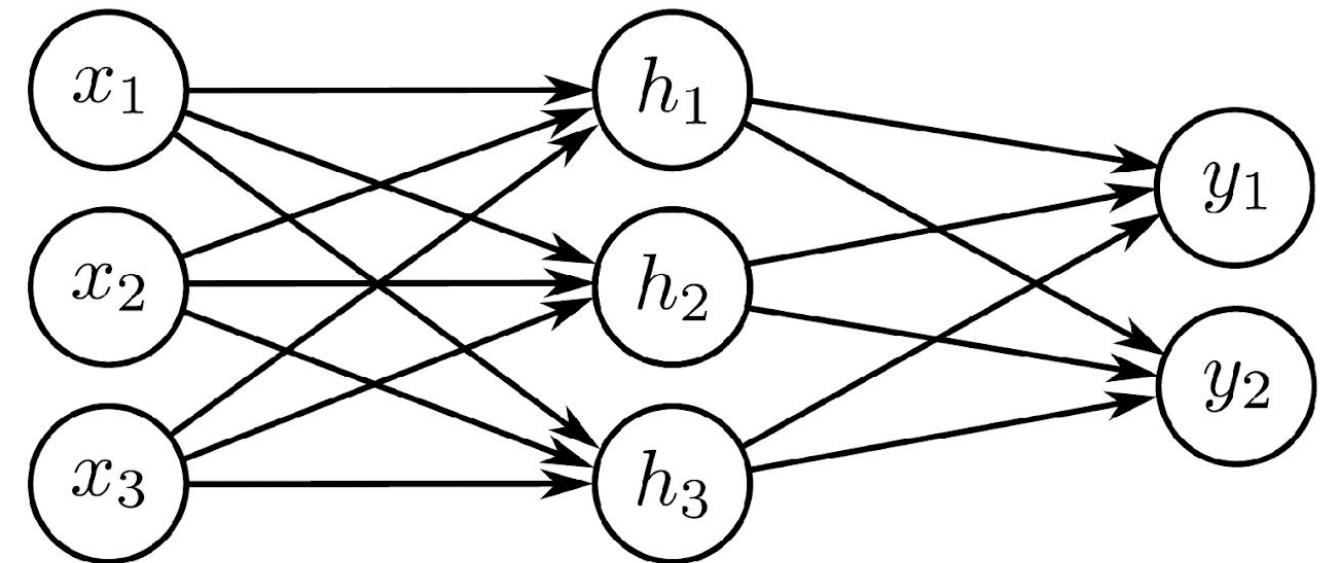
**Figure 3.8** Processing in network with two inputs  $\mathbf{x} = [x_1, x_2]^T$ , three hidden units  $h_1, h_2, h_3$ , and one output  $y$ . a–c) The input to each hidden unit is a linear function of the two inputs, which corresponds to an oriented plane. Brightness indicates function output. For example, in panel (a), the brightness represents  $\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2$ . Thin lines are contours. d–f) Each plane is clipped by the ReLU activation function (cyan lines are equivalent to “joints” in figures 3.3d–f). g–i) The clipped planes are then weighted, and j) summed together with an offset that determines the overall height of the surface. The result is a continuous surface made up of convex piecewise linear polygonal regions.



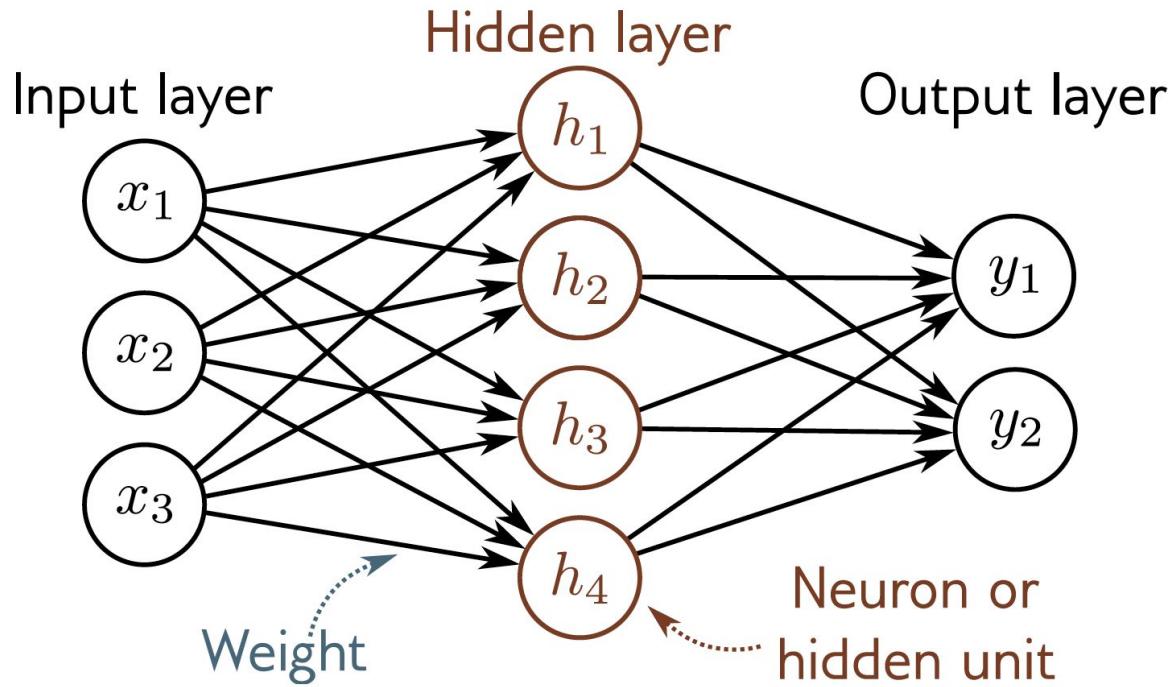
**Figure 3.9** Linear regions vs. hidden units. a) Maximum possible regions as a function of the number of hidden units for five different input dimensions  $D_i = \{1, 5, 10, 50, 100\}$ . The number of regions increases rapidly in high dimensions; with  $D = 500$  units and input size  $D_i = 100$ , there can be greater than  $10^{107}$  regions (solid circle). b) The same data are plotted as a function of the number of parameters. The solid circle represents the same model as in panel (a) with  $D = 500$  hidden units. This network has 51,001 parameters and would be considered very small by modern standards.



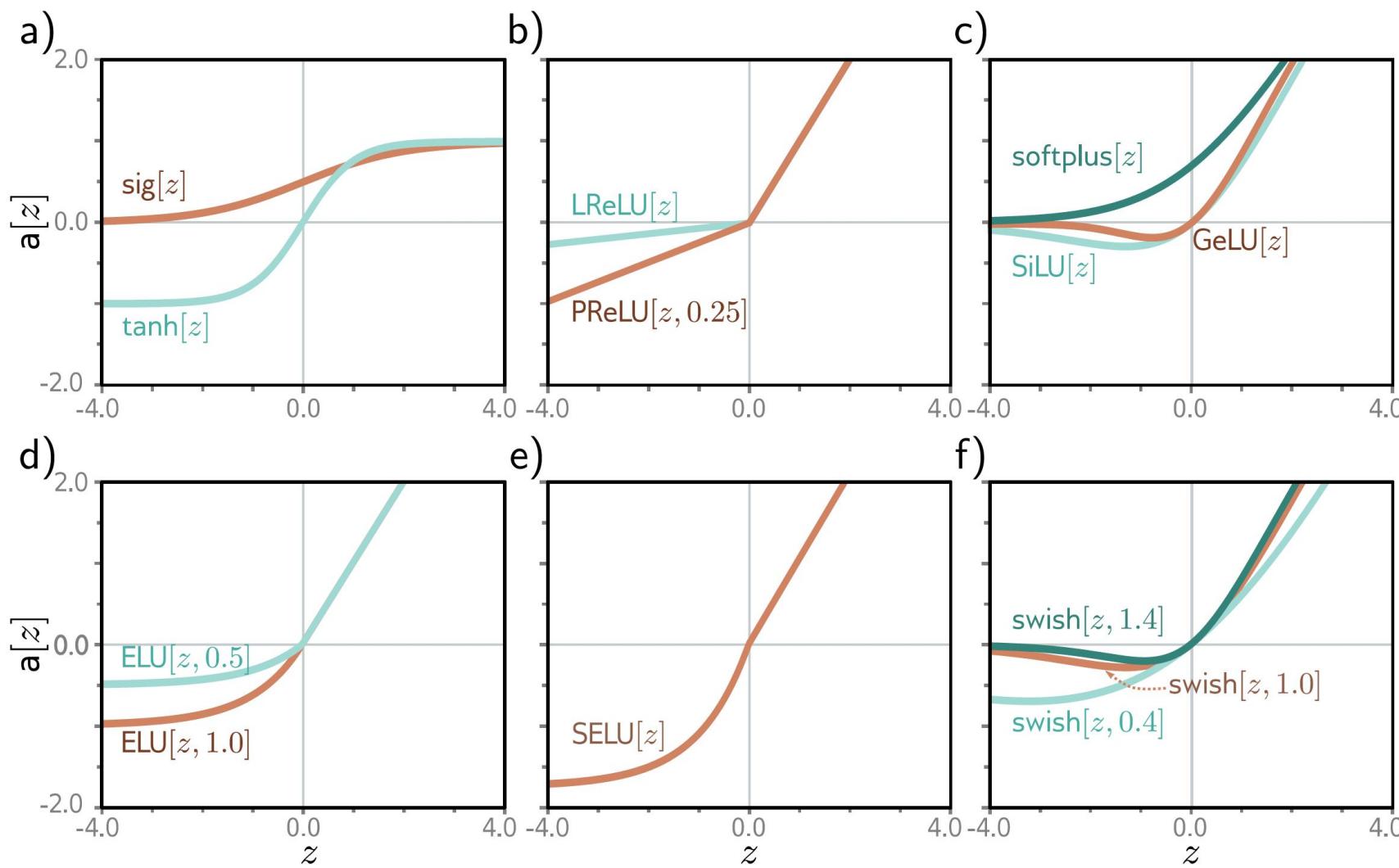
**Figure 3.10** Number of linear regions vs. input dimensions. a) With a single input dimension, a model with one hidden unit creates one joint, which divides the axis into two linear regions. b) With two input dimensions, a model with two hidden units can divide the input space using two lines (here aligned with axes) to create four regions. c) With three input dimensions, a model with three hidden units can divide the input space using three planes (again aligned with axes) to create eight regions. Continuing this argument, it follows that a model with  $D_i$  input dimensions and  $D_i$  hidden units can divide the input space with  $D_i$  hyperplanes to create  $2^{D_i}$  linear regions.



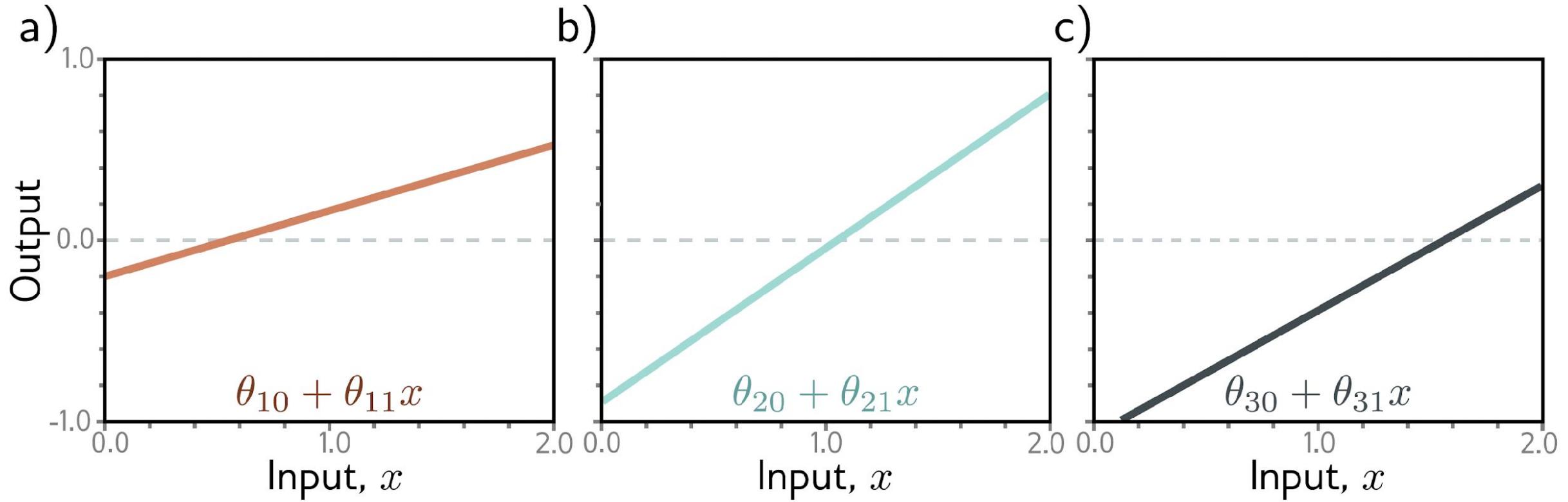
**Figure 3.11** Visualization of neural network with three inputs and two outputs. This network has twenty parameters. There are fifteen slopes (indicated by arrows) and five offsets (not shown).



**Figure 3.12** Terminology. A shallow network consists of an input layer, a hidden layer, and an output layer. Each layer is connected to the next by forward connections (arrows). For this reason, these models are referred to as feed-forward networks. When every variable in one layer connects to every variable in the next, we call this a fully connected network. Each connection represents a slope parameter in the underlying equation, and these parameters are termed weights. The variables in the hidden layer are termed neurons or hidden units. The values feeding into the hidden units are termed pre-activations, and the values at the hidden units (i.e., after the ReLU function is applied) are termed activations.



**Figure 3.13** Activation functions. a) Logistic sigmoid and tanh functions. b) Leaky ReLU and parametric ReLU with parameter 0.25. c) SoftPlus, Gaussian error linear unit, and sigmoid linear unit. d) Exponential linear unit with parameters 0.5 and 1.0, e) Scaled exponential linear unit. f) Swish with parameters 0.4, 1.0, and 1.4.



**Figure 3.14** Processing in network with one input, three hidden units, and one output for problem 3.4. a–c) The input to each hidden unit is a linear function of the inputs. The first two are the same as in figure 3.3, but the last one differs.