**2.**



**Figure 4.6** Matrix notation for network with $D_i = 3$-dimensional input $\mathbf{x}$, $D_o = 2$-dimensional output $\mathbf{y}$, and $K = 3$ hidden layers $\mathbf{h}_1, \mathbf{h}_2,$ and $\mathbf{h}_3$ of dimensions $D_1 = 4$, $D_2 = 2$, and $D_3 = 3$ respectively. The weights are stored in matrices $\mathbf{\Omega}_k$ that multiply the activations from the preceding layer to create the pre-activations at the subsequent layer. For example, the weight matrix $\mathbf{\Omega}_1$ that computes the pre-activations at $\mathbf{h}_2$ from the activations at $\mathbf{h}_1$ has dimension $2 \times 4$. It is applied to the four hidden units in layer one and creates the inputs to the two hidden units at layer two. The biases are stored in vectors $\boldsymbol{\beta}_k$ and have the dimension of the layer into which they feed. For example, the bias vector $\boldsymbol{\beta}_2$ is length three because layer $\mathbf{h}_3$ contains three hidden units.

- Number of layers
- Hidden units in layer 1
- Hidden units in layer 2
- Hidden units in layer 3

**4.**



$$\underline{h}_1 = a\left[\underline{\beta}_0 + \underline{\Omega}_0 \, \underline{x}\right]$$

$$\underline{h}_2 = a\left[\underline{\beta}_1 + \underline{\Omega}_1 \, \underline{h}_1\right]$$

$$\underline{h}_3 = a\left[\underline{\beta}_2 + \underline{\Omega}_2 \, \underline{h}_2\right]$$

$$\underline{y} = \underset{\uparrow}{O}\left[\underline{\beta}_3 + \underline{\Omega}_3 \, \underline{h}_3\right]$$

some output activation func.

$$x \rightarrow 5 \times 1$$
$$\underline{\Omega}_0 \rightarrow 20 \times 5$$
$$\beta_0 \rightarrow 20 \times 1$$
$$\underline{\Omega}_1 \rightarrow 10 \times 20$$
$$\beta_1 \rightarrow 10 \times 1$$
$$\underline{\Omega}_2 \rightarrow 7 \times 10$$
$$\beta_2 \rightarrow 7 \times 1$$
$$\underline{\Omega}_3 \rightarrow 4 \times 7$$
$$\beta_3 \rightarrow 4 \times 1$$

6. ( consider just weights, not biases).

There are ccurrently,

$1 \times 10 + 9 (10 \times 10) + 1 \times 10 = 10 + 900 + 10 = 920$ weights
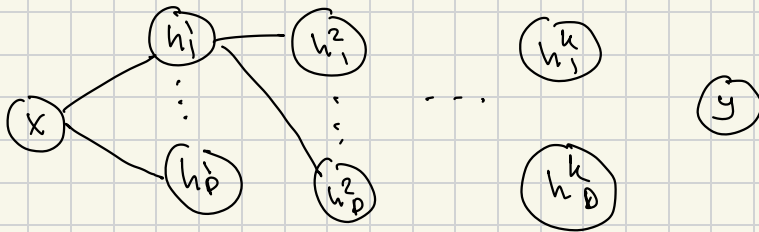Increasing a layer (depth) gives 100 more weights $= 1020$ weights
Increasing width (no nodes in each layer) gives

$1 \times 11 + 9 (11 \times 11) + 1 \times 11 = 1111$ weights.

So, more weights added by increasing width.

8. In the 'odd' answers.

10. Looks like.



Between the input and layer 1 there are $D + D$ params.

Between final layer and output there are $D + 1$ params

In between layers there are $(k-1)(D \times D + D)$ params.

this becomes

$D + D + D + 1 + (k-1)(D \times D + D)$
$= 3D + 1 + (k-1) D (D+1)$ parameters.