Deep MLP training: stochastic gradient descent.

**DEEP-MLP-TRAINING ($\mathbf{D}, h, \eta, \texttt{maxiter}, n_1, n_2, \cdots, n_h, f^1, f^2, \cdots, f^{h+1}$):**

1   $n_0 \leftarrow d$ // input layer size
2   $n_{h+1} \leftarrow p$ // output layer size
    // Initialize weight matrices and bias vectors
3   **for** $l = 0, 1, 2, \cdots, h$ **do**
4     |   $\boldsymbol{\theta}_l \leftarrow$ random $n_{l+1}$ vector with small values
5     |   $\mathbf{W}_l \leftarrow$ random $n_l \times n_{l+1}$ matrix with small values
6   $t \leftarrow 0$ // iteration counter
7   **repeat**
8     |   **foreach** $(\mathbf{x}_i, \mathbf{y}_i) \in \mathbf{D}$ *in random order* **do**
          // Feed-Forward Phase
9         |   $\mathbf{z}^0 \leftarrow \mathbf{x}_i$
10       |   **for** $l = 0, 1, 2, \ldots, h$ **do**
11         |   $\mathbf{z}^{l+1} \leftarrow f^{l+1}\left(\mathbf{W}_l^T \cdot \mathbf{z}^l\right)$
12       |   $\mathbf{o}_i \leftarrow \mathbf{z}^{h+1}$
          // Backpropagation Phase
13       |   $\boldsymbol{\delta}^{h+1} \leftarrow \partial \mathbf{f}^{h+1} \odot \partial \boldsymbol{\mathcal{E}}_{\mathbf{x}}$ // net gradients at output
14                 // use $\partial \mathbf{F}^{h+1} \partial \boldsymbol{\mathcal{E}}_{\mathbf{x}}$ for softmax
15       |   **for** $l = h, h-1, \cdots, 1$ **do**
16         |   $\boldsymbol{\delta}^l \leftarrow \partial \mathbf{f}^l \odot \left(\mathbf{W}_l \cdot \boldsymbol{\delta}^{l+1}\right)$ // net gradients at layer $l$
          // Gradient Descent Step
17       |   **for** $l = 0, 1, \cdots, h$ **do**
18         |   $\nabla_{\mathbf{W}_l} \leftarrow \mathbf{z}^l \cdot \left(\boldsymbol{\delta}^{l+1}\right)^T$ // weight gradient matrix at layer $l$
19         |   $\nabla_{\boldsymbol{\theta}_l} \leftarrow \boldsymbol{\delta}^{l+1}$ // bias gradient vector at layer $l$
20       |   **for** $l = 0, 1, \cdots, h$ **do**
21         |   $\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \cdot \nabla_{\mathbf{W}_l}$ // update $\mathbf{W}_l$
22         |   $\boldsymbol{\theta}_l \leftarrow \boldsymbol{\theta}_l - \eta \cdot \nabla_{\boldsymbol{\theta}_l}$ // update $\boldsymbol{\theta}_l$
23     |   $t \leftarrow t + 1$
24   **until** $t \geq \texttt{maxiter}$