**Input:** loss function $E$, learning rate $\eta$, dataset $X, y$ und das Modell
$\quad\quad$ $F(\theta, x)$
**Output:** Optimum $\theta$ which minimizes $\epsilon$

**1** **while** *converge* **do**
**2** $\quad$ $\tilde{y} = F(\theta, x)$
**3** $\quad$ $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^{N} \frac{\delta\epsilon(y, \tilde{y})}{\delta\theta}$
**4** **end**

**Algorithm 1:** Gradient descent

**Input:** loss function $E$, learning rate $\eta$, dataset $X, y$ und das Modell
$\quad\quad$ $F(\theta, x)$
**Output:** Optimum $\theta$ which minimizes $\epsilon$

**1** **while** *converge* **do**
**2** $\quad$ Shuffle X, y
**3** $\quad$ **for** $x_i, y_i$ *in X, y* **do**
**4** $\quad\quad$ $\tilde{y} = F(\theta, x_i)$
**5** $\quad\quad$ $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^{N} \frac{\delta\epsilon(y_i, \tilde{y_i})}{\delta\theta}$
**6** $\quad$ **end**
**7** **end**

**Algorithm 2:** Stochastic Gradient descent(SGD)

**Input:** loss function $E$, learning rate $\eta$, dataset $X, y$ und das Modell
$\quad\quad F(\theta, x)$
**Output:** Optimum $\theta$ which minimizes $E$
**1 while** *converge* **do**
**2** $\quad$ Shuffle X, y
**3** $\quad$ **for** *each batch of $x_i, y_i$ in X, y* **do**
**4** $\quad\quad$ $\tilde{y} = F(\theta, x_i)$
**5** $\quad\quad$ $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^{N} \frac{\delta E(y_i, \tilde{y_i})}{\delta E}$
**6** $\quad$ **end**
**7 end**

**Algorithm 3:** Mini-Batch Stochastic Gradient descent(MSGD)

**Input:** Netzwerk mit $l$ layers, Aktivirungsfunktion $\sigma_l$ , Output von der
$\quad\quad$ verstekten Schicht $h_l = \sigma_l(W_l^T h_{l-1} + b_l)$ und die
$\quad\quad$ Netzwerkausgabe $\tilde{y} = h_l$
**1** Berechnen der Gradient: $\delta \leftarrow \frac{\partial E(y_i, \tilde{y_i})}{\partial y}$
**2 for** $i \leftarrow l$ *bis* $0$ **do**
**3** $\quad$ Berechnen der Gradient für die Aktuelle Schicht
**4** $\quad$ $\frac{\partial E(y,\tilde{y})}{\partial W_l} = \frac{\partial E(y,\tilde{y})}{\partial h_l} \frac{\partial h_l}{\partial W_l} = \delta \frac{\partial h_l}{\partial W_l}$
**5** $\quad$ $\frac{\partial E(y,\tilde{y})}{\partial b_l} = \frac{\partial E(y,\tilde{y})}{\partial h_l} \frac{\partial h_l}{\partial b_l} = \delta \frac{\partial h_l}{\partial b_l}$
**6** $\quad$ Gradientabstiegverfahren mit $\frac{\partial E(y,\tilde{y})}{\partial W_l}$ und $\frac{\partial E(y,\tilde{y})}{\partial b_l}$
**7** $\quad$ Propagiere den Gradienten zu den unteren Schichten.
**8** $\quad$ $\delta \leftarrow \frac{\partial E(y,\tilde{y})}{\partial h_l} \frac{\partial h_l}{\partial h_{l-1}} = \delta \frac{\partial h_l}{\partial h_{l-1}}$
**9 end**

**Algorithm 4:** Back-Propagation