



CS F425: Deep Learning

10

Backpropagation Training



Dr. Kamlesh Tiwari
 Assistant Professor, Department of CSIS,
 BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
 Feb 10, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

Recap: Backpropagation (for 2 layers)

Algorithm 1: Backpropagation($D, \eta, \eta_{in}, \eta_{out}, \eta_{hidden}$)

- 1 Create the feedforward network with $\eta_{in}, \eta_{out}, \eta_{hidden}$ layers
 - 2 Randomly initialize weights to small values $\in [-0.05, +0.05]$
 - 3 **repeat**
 - 4 **for each** $\langle \vec{x}, \vec{t} \rangle \in D$ **do**
 - 5 $o_u =$ get output from network \forall unit u
 - 6 $\delta_k = o_k(1 - o_k)(t_k - o_k)$ for all **output unit** k
 - 7 $\delta_h = o_h(1 - o_h) \sum_{k \in \text{outputs}} (W_{kh} \delta_k)$ for all **hidden unit** h
 - 8 $w_{ji} = w_{ji} + \Delta w_{ji}$ where $\Delta w_{ji} = \eta \delta_j x_{ji}$
 - 9 **until converge;**
- Recall error function is $E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$
 - For a single training example $E_d(\vec{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$
 - Weight w_{ji} is updated by adding $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$

Value of $\frac{\partial E_d}{\partial net_j}$ for output units

- $\frac{\partial E_d}{\partial net_j} = \frac{\partial o_j}{\partial \sigma_j} \times \frac{\partial o_j}{\partial net_j}$
 - $\frac{\partial E_d}{\partial \sigma_j} = \frac{\partial}{\partial \sigma_j} \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2 = \frac{\partial}{\partial \sigma_j} \frac{1}{2} (t_j - o_j)^2 = -(t_j - o_j)$
 - Note that $o_j = \sigma(net_j)$ therefore $\frac{\partial o_j}{\partial net_j}$ is derivative of sigmoid
- $$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \frac{1}{1 + e^{-x}} = \frac{d}{dx} (1 + e^{-x})^{-2} \frac{d}{dx} (1 + e^{-x}) \\ &= (-1)(1 + e^{-x})^{-2} (0 - e^{-x}) \\ &= \frac{1}{1 + e^{-x}} \times \frac{e^{-x}}{1 + e^{-x}} = \sigma(x)(1 - \sigma(x)) \end{aligned}$$
- As a result $\frac{\partial o_j}{\partial net_j} = \frac{\partial \sigma(net_j)}{\partial net_j} = \sigma(net_j)(1 - \sigma(net_j)) = o_j(1 - o_j)$
 - $\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j(1 - o_j)$
 - Term $(t_j - o_j) o_j(1 - o_j)$ is treated as δ_j

Therefore, $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} \times x_{ji} = \eta(t_j - o_j) o_j(1 - o_j) x_{ji}$

Recap: Training Multilayer Networks with Backpropagation

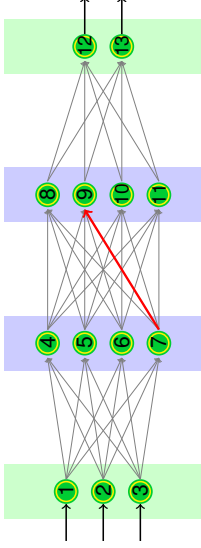
- Single perceptron can only express linear decision surface
 - We need units whose output is a nonlinear function of its inputs AND is also differentiable (using Neuron not Perceptron)
- $o(\vec{x}) = \sigma(\vec{w} \cdot \vec{x})$ where $\sigma(y) = \frac{1}{1 + e^{-y}}$

Backpropagation algorithm learns weights for a fixed set of units and interconnections

- It employs **gradient descent** to minimize the error between the network output values and the target values for these outputs
- Let Error function is redefined as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$$

Recap: Conventions Over The Network



- x_{ji} i th input to unit j (x_{44} is highlighted)
- w_{ji} weight associated with i th input to unit j
- net_j be $\sum_i w_{ji} x_{ji}$ the weighted sum of input for unit j
- o_j output computed by unit j . Consider it as $\sigma(net_j)$
- t_j target output for unit j
- outputs** set of units in final layer ($\{12, 13\}$ in our case)
- Downstream(l)** units whose immediate input is the output of unit j

We are interested in $\frac{\partial E_d}{\partial w_{ji}}$ it is $\frac{\partial E_d}{\partial net_j} \times \frac{\partial net_j}{\partial w_{ji}}$ and therefore, $\frac{\partial E_d}{\partial net_j}$

Value of $\frac{\partial E_d}{\partial net_j}$ for hidden units

$$\begin{aligned} \frac{\partial E_d}{\partial net_j} &= \sum_{k \in \text{Downstream}(j)} \frac{\partial E_d}{\partial net_k} \times \frac{\partial net_k}{\partial net_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times \frac{\partial net_k}{\partial net_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times \frac{\partial net_k}{\partial \sigma_j} \times \frac{\partial \sigma_j}{\partial net_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times w_{kj} \times o_j(1 - o_j) \end{aligned}$$

- δ_j being $-\frac{\partial E_d}{\partial net_j} = o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k \times w_{kj}$

Therefore, $\Delta w_{ji} = \eta \delta_j x_{ji} = \eta(o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k \times w_{kj}) x_{ji}$

Note: $net_j = w_{j1}x_{j1} + w_{j2}x_{j2} + \dots + w_{jn}x_{jn}$ so o_j is input to k th unit so $o_j = x_{kj}$ so $\frac{\partial net_k}{\partial net_j} = w_{kj}$

The Backpropagation Algorithm

Algorithm 2: Backpropagation ($D, \eta, n_{in}, n_{out}, n_{hidden}$)

- 1 Create the feedforward network with $n_{in}, n_{out}, n_{hidden}$ layers
- 2 Randomly initialize weights to small values $\in [-0.05, +0.05]$
- 3 **repeat**
- 4 **for each** $\langle \vec{x}, \vec{t} \rangle \in D$ **do**
- 5 o_u = get output from network \forall unit u
- 6 $\delta_k = o_k(1 - o_k)(t_k - o_k)$ for all **output unit** k
- 7 $\delta_h = o_h(1 - o_h) \sum_{k \in \text{outputs}} (w_{hk} \delta_k)$ for all **hidden unit** h
- 8 $w_{ji} = w_{ji} + \Delta w_{ji}$ where $\Delta w_{ji} = \eta \delta_j x_{ji}$
- 9 **until converge;**

Backpropagation

- **Adding Momentum:** weight update during n^{th} iteration depend partially on the update that occurred during the $(n - 1)^{\text{th}}$ iteration

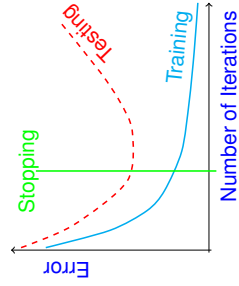
$$\Delta w_{ji}(n) = \eta \delta_j x_{ji} + \alpha \Delta w_{ji}(n-1)$$

- **Learning in arbitrary acyclic network:** for feed-forward networks of arbitrary depth, δ_r value for a unit in hidden layer is determined as

$$\delta_r = o_r(1 - o_r) \sum_{s \in \text{Downstream}(r)} w_{sr} \times \delta_s$$

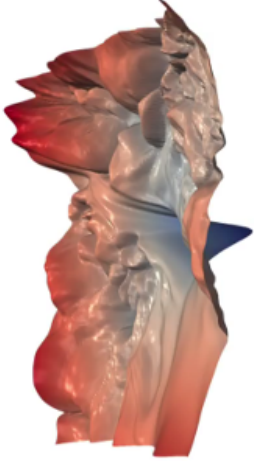
Generalization, Overfitting, and Stopping Criterion

Continue training until the error on the training examples falls below some predetermined threshold could be a poor strategy



- Weight decay or use of validation set (k-fold ?) is suggested
- Input or output encoding can be used

Loss Landscape 1



Backpropagation over multilayer networks converge to a local minimum, **NOT necessarily to the global minima**

Backpropagation

- Result of Backpropagation over multilayer networks is only guaranteed to converge toward some local minimum and not necessarily to the global minimum error
- No methods are known to predict with certainty when local minima will cause difficulties
- Suggested to use momentum, true gradient descent or multiple networks (initialized with different random weights)
- Any boolean function can precisely be represented by some network having only **two** layers of units (Cybenko 1989; Hornik et al. 1989)
- Every bounded continuous function can be approximated with arbitrarily small error (under a finite norm) by a network with two layers of units
- Any function can be approximated to arbitrary accuracy by a network with three layers of units (Cybenko 1988)

Coding Example 2

Consider two class data
Logistic Regression
Neural Network
Decision Boundary
Bigger hidden layer

