



UMassAmherst
The Commonwealth's Flagship Campus

Neural Networks

Chaitra Gopalappa

Backprop (back propagation)

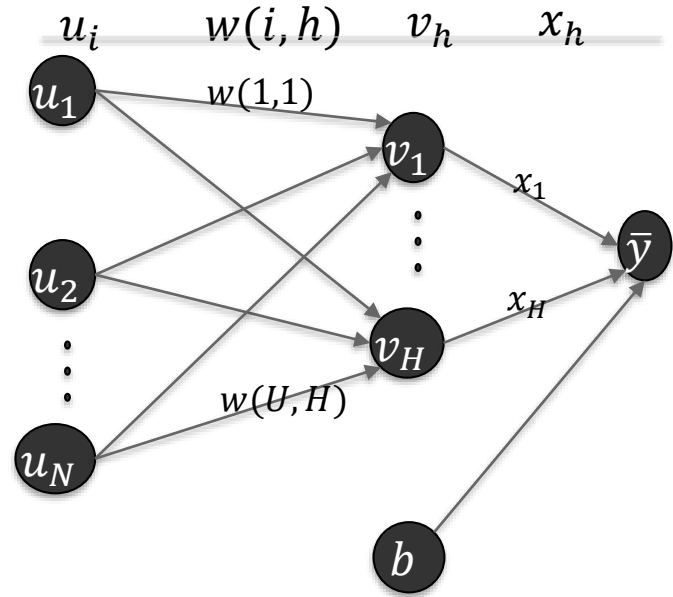
Example with one hidden layer and sigmoid activation

- Backprop concepts

- Objective

$$\min_{b,x,w} L = \min_{b,x,w} \frac{1}{2} \sum_{p=1}^P (y_p - \bar{y}_p)^2$$

- $b, \omega[.,.],$ and $x[.]$ are the decision variables from perspective of Backprop.
- \bar{y}_p = output from trained net for \vec{u}_p ; p are the data samples
- y_p = actual data
- $L = f(b, \omega[.,.], x[.])$ is the lost function
- P = number of samples; N = number of input nodes



Backprop algorithm

1. Initialize

1. Initialize b , $\omega(\cdot, \cdot)$, and $x(\cdot)$ to random values
2. Set SSE_{old} to very large value.
3. Set $m = 0$ (iteration number)

2. Forward pass

1. Compute $v_p^*[h] = \sum_i w[i, h]u_p[i]$ for each p and each h .
2. Compute $v_p[h] = \frac{1}{1+e^{-v_p^*[h]}}$ for each p and each h
3. Compute $\bar{y}_p = b + \sum_h x[h]v_p[h]$ for each p (data samples)

3. Apply SD transformations

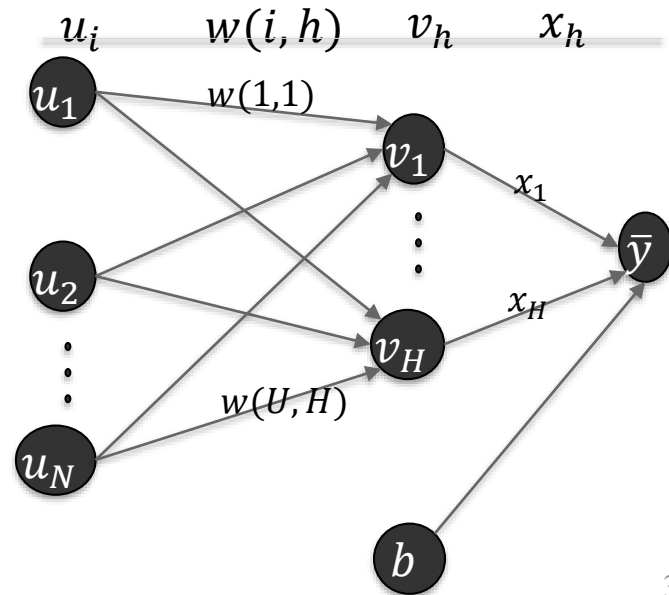
$$b_m \leftarrow b_{m-1} - \mu \frac{\partial(L)}{\partial b}$$
$$\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \mu \frac{\partial(L)}{\partial \omega[i, h]}$$
$$x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(L)}{\partial x[h]}$$

4. Set $m = m + 1$

1. Calculate $L_{new} = \sum_p (y_p - \bar{y}_p)^2$
2. Update μ
3. If $[L_{new} - L] < \text{tolerance} \rightarrow STOP$ Otherwise set $L_{old} = L_{new}$ got step 2.

Derivation of Backprop derivatives

- Main transformation is the steepest descent $\vec{x} \leftarrow \vec{x} - \mu \nabla f(\vec{x})$; but here we are solving for $\omega[.,.]$, and $x[.]$
- $b_m \leftarrow b_{m-1} - \mu \frac{\partial(L)}{\partial b}$
- $\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \frac{\mu \partial(L)}{\partial \omega[i, h]}$
- $x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(L)}{\partial x[h]}$



$\frac{\partial L}{\partial x[h]}=? \frac{\partial L}{\partial b}=?$ **Apply chain rule of calculus**

$$\begin{aligned} \bullet \quad \frac{\partial L}{\partial x[h]} &= \frac{1}{2} \sum_{p=1}^P \frac{\partial}{\partial x[h]} (y_p - \bar{y}_p)^2 \\ &= \frac{2}{2} \sum_{p=1}^P (y_p - \bar{y}_p) \frac{\partial}{\partial x[h]} (y_p - \bar{y}_p) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial \bar{y}_p}{\partial x[h]} \right) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial (b + \sum_h x[h] v_p[h])}{\partial x[h]} \right) \end{aligned}$$

$$\frac{\partial L}{\partial x[h]} = -\sum_p (y_p - \bar{y}_p) v_p[h]$$

$$\bullet \quad \frac{\partial L}{\partial b} = -\sum_p (y_p - \bar{y}_p)$$

$$\frac{\partial L}{\partial \omega(i, h)} = ?$$

$$\begin{aligned} \frac{\partial L}{\partial \omega[i, h]} &= \frac{1}{2} \sum_p \frac{\partial}{\partial \omega[i, h]} (y_p - \bar{y}_p)^2 \\ &= \frac{1}{2} \sum_p 2(y_p - \bar{y}_p) \frac{\partial}{\partial \omega[i, h]} (y_p - \bar{y}_p) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial}{\partial w[i, h]} \bar{y}_p \right) \\ &= - \sum_p (y_p - \bar{y}_p) \cdot \left(\frac{\partial \bar{y}_p}{\partial v_p[h]} \frac{\partial v_p[h]}{\partial v_p^*[h]} \frac{\partial v_p^*[h]}{\partial \omega[i, h]} \right) \\ \frac{\partial L}{\partial w[i, h]} &= - \sum_p ((y_p - \bar{y}_p) x[h] v_p[h] (1 - v_p[h]) u_p[i]) \end{aligned}$$

This changes for
different activation
functions

- Feed forward equations
 - $v_p^*[h] = \sum_i \omega[i, h] u_p[i];$
 - $v_p[h] = \frac{1}{1 + e^{-v_p^*[h]}};$
 - $y_p = b + \sum_h v_p[h] x[h]$
- $\frac{\partial v_p^*[h]}{\partial w(i, h)} = u_p[i]$
- $\frac{\partial v_p[h]}{\partial v_p^*[h]} = \frac{\partial (1 + e^{-v_p^*[h]})^{-1}}{\partial v_p^*[h]}$

$$= \frac{-1}{[1 + e^{-v_p^*[h]}]^2} e^{-v_p^*[h]} (-1)$$

$$= \frac{1 + e^{-v_p^*[h]} - 1}{[1 + e^{-v_p^*[h]}]^2} \text{ (add + 1, - 1 in numerator)}$$

$$= \frac{1 + e^{-v_p^*[h]}}{[1 + e^{-v_p^*[h]}]^2} - \frac{1}{[1 + e^{-v_p^*[h]}]^2}$$

$$= \frac{1}{1 + e^{-v_p^*[h]}} - \frac{1}{[1 + e^{-v_p^*[h]}]^2}$$

$$= v_p[h] (1 - v_p[h])$$

Backprop algorithm

1. Initialize

1. Initialize $b, \omega(\cdot, \cdot)$, and $x(\cdot)$ to random values
2. Set SSE_{old} to very large value.
3. Set $m = 0$ (iteration number)

2. Forward pass

1. Compute $v_p^*[h] = \sum_i w[i, h]u_p[i]$ for each p and each h .
2. Compute $v_p[h] = \frac{1}{1+e^{-v_p^*[h]}}$ for each p and each h
3. Compute $\bar{y}_p = b + \sum_h x[h]v_p[h]$ for each p (data samples)

3. Apply SD transformations

$$\begin{aligned}
 b_m &\leftarrow b_{m-1} - \mu \frac{\partial(L)}{\partial b}; & \frac{\partial L}{\partial b} &= -\sum_p (y_p - \bar{y}_p) \\
 \omega_m[i, h] &\leftarrow \omega_{m-1}[i, h] - \mu \frac{\partial(L)}{\partial \omega[i, h]}; & \frac{\partial L}{\partial \omega[i, h]} &= -\sum_p ((y_p - \bar{y}_p)x[h]v_p[h](1 - v_p[h])u_p[i]) \\
 x_m[h] &\leftarrow x_{m-1}[h] - \mu \frac{\partial(L)}{\partial x[h]}; & \frac{\partial L}{\partial x[h]} &= -\sum_p (y_p - \bar{y}_p)v_p[h]
 \end{aligned}$$

4. Set $m = m + 1$

1. Calculate $L_{new} = \sum_p (y_p - \bar{y}_p)^2$
2. Update μ
3. If $[L_{new} - L_{old}] < \text{tolerance} \rightarrow STOP$ Otherwise set $L_{old} = L_{new}$ got step 2.

- Do multiple iterations of the algorithm to start at different initial points
- Initialize $b, \omega[\cdot, \cdot]$, and $x[\cdot]$ to select from a random range; Try a few different ranges
- Try a few different options for learning rate ($\mu = \text{small constant}; \mu = \frac{A}{B+m}; \mu = \frac{1}{m};$)

This applies to sigmoid activation; change if using a different activation function

UMassAmherst
The Commonwealth's Flagship Campus