



CS F425: Deep Learning

09

Delta Training Rule Activation Function



Dr. Kamlesh Tiwari
Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Feb 09, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

Perceptron Training (delta rule)

Algorithm 1: Gradient Descent (D, η)

- 1 Initialize w_i with random weights
- 2 **repeat**
- 3 For each w_i , initialize $\Delta w_i = 0$
- 4 **for** each training example $d \in D$ **do**
- 5 Compute output o using model for d whose target is t
- 6 For each w_i , update $\Delta w_i = \Delta w_i + \eta(t-o)x_i$
- 7 For each w_i , set $w_i = w_i + \Delta w_i$
- 8 **until** termination condition is met;
- 9 **return** w

- A date item $d \in D$, is supposed to be multidimensional $d = (x_1, x_2, \dots, x_n, t)$
- Algorithm converges toward the minimum error hypothesis.
- Linear programming can also be an approach

Perceptron Training (delta rule)

When data is not linearly-separable; error fluctuates with parameter training updates. It is difficult to decide, when to stop?

- **Delta rule** converges to a best-fit approximation of the target
- Uses **gradient descent**
- Consider unthresholded perceptron, $o(\vec{x}) = \vec{w} \cdot \vec{x}$
- Training error is defined as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- Gradient would specify direction of steepest increase $\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$
- Weights can be learned as $w_j = w_j - \eta \frac{\partial E}{\partial w_j}$
- It can be seen that $\frac{\partial E}{\partial w_j} = \sum_{d \in D} (t_d - o_d)(-x_{jd})$

Perceptron Training (delta rule)

When data is not linearly-separable; error fluctuates with parameter training updates. It is difficult to decide, when to stop?

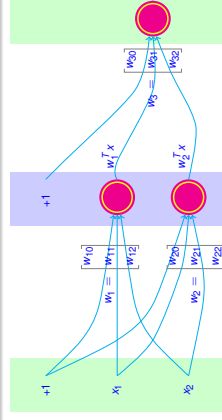
- **Delta rule** converges to a best-fit approximation of the target
- Uses **gradient descent**
- Consider unthresholded perceptron, $o(\vec{x}) = \vec{w} \cdot \vec{x}$
- Training error is defined as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- Gradient would specify direction of steepest increase $\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$
- Weights can be learned as $w_j = w_j - \eta \frac{\partial E}{\partial w_j}$
- It can be seen that $\frac{\partial E}{\partial w_j} = \sum_{d \in D} (t_d - o_d)(-x_{jd})$

Linear Activation is Not Much Interesting

NN with perceptrons have limited capability, even with many layers



$$\begin{aligned} \text{Output} &= w_{00} \times 1 + w_{01} \times (w_1^T x) + w_{02} \times (w_2^T x) \\ &= w_{00} \times 1 + w_{01} \times [w_{10} \times 1 + w_{11} \times x_1 + w_{12} \times x_2] \\ &\quad + w_{02} \times [w_{20} \times 1 + w_{21} \times x_1 + w_{22} \times x_2] \\ &= (w_{00} + w_{01}w_{10} + w_{02}w_{20}) + (w_{01}w_{11} + w_{02}w_{21}) \times x_1 \\ &\quad + (w_{01}w_{12} + w_{02}w_{22}) \times x_2 \\ &= w_0^T \times 1 + w_1^T \times x_1 + w_2^T \times x_2 \end{aligned}$$

Expression of single perceptron

Perceptron Training (delta rule)

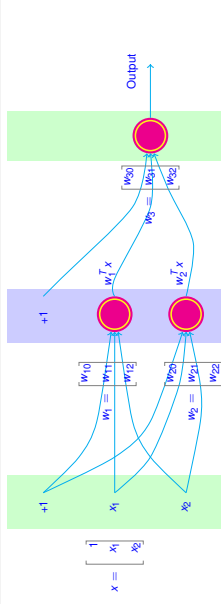
Algorithm 2: Gradient Descent (D, η)

- 1 Initialize w_i with random weights
- 2 **repeat**
- 3 For each w_i , initialize $\Delta w_i = 0$
- 4 **for** each training example $d \in D$ **do**
- 5 Compute output o using model for d whose target is t
- 6 For each w_i , update $\Delta w_i = \Delta w_i + \eta(t-o)x_i$
- 7 For each w_i , set $w_i = w_i + \Delta w_i$
- 8 **until** termination condition is met;
- 9 **return** w

- A date item $d \in D$, is supposed to be multidimensional $d = (x_1, x_2, \dots, x_n, t)$
- Algorithm converges toward the minimum error hypothesis.
- Linear programming can also be an approach

Limitation of Linear Activation

NN with perceptrons have limited capability, even with many layers

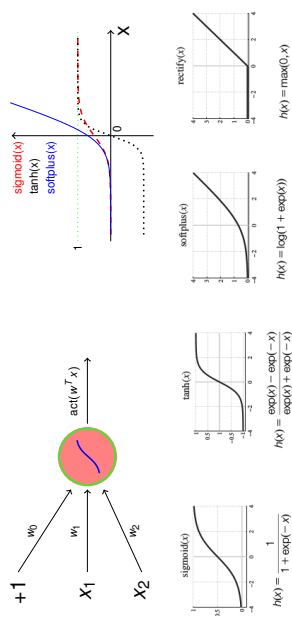


$$\begin{aligned} \text{Output} &= w_{00} \times 1 + w_{01} \times (w_{10} \times 1 + w_{12} \times (w_2^T \cdot x)) \\ &= w_{00} \times 1 + w_{01} \times (w_{10} \times 1 + w_{11} \times x_1 + w_{12} \times x_2) \\ &= (w_{00} + w_{01} w_{10} + w_{02} w_{20}) + (w_{01} w_{11} + w_{02} w_{21}) \times x_1 \\ &\quad + (w_{01} w_{12} + w_{02} w_{22}) \times x_2 \\ &= w_0^T \times x_1 + w_2^T \times x_2 \end{aligned}$$

Expression of single perceptron

Neuron

Neuron uses nonlinear **activation functions** (*sigmoid*, *tanh*, *softplus*, *ReLU* leaky *ReLU* etc.) at the place of thresholding



To avoid dying ReLU problem, leakyReLU is used $\max(0, 1x, \lambda x)$

More Activation Functions

ReLU $\max(0, x)$	GELU $\frac{x}{2} \left(1 + \tanh\left(\frac{\beta}{2} \left(x + \frac{2x^3}{5} + \dots\right)\right)\right)$	PReLU $\max(0, x)$	Soft Shrink $\begin{cases} x - \lambda & \text{if } x > \lambda \\ x + \lambda & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$
ELU $\begin{cases} x & \text{if } x > 0 \\ \alpha \exp(x - 1) & \text{if } x < 0 \end{cases}$	Swish $\frac{x}{1 + \exp(-x)}$	SELU $\alpha \max(0, x) + \min(0, \beta(\exp(x - 1)))$	Hard Shrink $\begin{cases} x & \text{if } x > \lambda \\ 0 & \text{otherwise} \end{cases}$
SoftPlus $\frac{1}{2} \log(1 + \exp(2x))$	Mish $x \tanh\left(\frac{\pi}{2} \log(1 + \exp(\pi x))\right)$	RReLU $\begin{cases} x & \text{if } x > 0 \\ \alpha & \text{if } x < 0 \text{ with } \alpha \sim \mathcal{U}(a, b) \end{cases}$	Soft Shrink $\begin{cases} x - \lambda & \text{if } x > \lambda \\ x + \lambda & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$
HardSigmoid $\begin{cases} 0 & \text{if } x \leq -3 \\ x & \text{if } x \in [-3, 3] \\ 1 & \text{if } x > 3 \end{cases}$	Sigmoid $\frac{1}{1 + \exp(-x)}$	Hard Sigmoid $\begin{cases} 0 & \text{if } x \leq -1 \\ x & \text{if } x \in [-1, 1] \\ 1 & \text{if } x > 1 \end{cases}$	Hard Shrink $\begin{cases} x & \text{if } x > \lambda \\ 0 & \text{otherwise} \end{cases}$
Tanh $\tanh(x)$	Hard tanh $\begin{cases} 1 & \text{if } x \geq a \\ x & \text{if } x \in [a, b] \\ -1 & \text{otherwise} \end{cases}$	Hard Sigmoid $\frac{x}{1 + x }$	Tanh Shrink $x - \tanh(x)$

Backpropagation (for 2 layers)

Algorithm 3: Backpropagation(D, η , η_{in} , η_{out} , η_{hidden})

- 1 Create the feedforward network with η_{in} , η_{out} , η_{hidden} layers
- 2 Randomly initialize weights to small values $\in [-0.05, +0.05]$
- 3 repeat
- 4 for each $\langle x, \tilde{y} \rangle \in D$ do
- 5 $O_j =$ 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 $w_{ij} = w_{ij} + \Delta w_{ij}$ where $\Delta w_{ij} = \eta \delta_j x_{ij}$
- 9 until *converge*;

- Recall error function is $E(\tilde{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$
- For a single training example $E_d(\tilde{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$
- Weight w_{ij} is updated by adding $\Delta w_{ij} = -\eta \delta_j x_{ij}$

Multilayer Networks and 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 (Use Neuron not Perceptron)

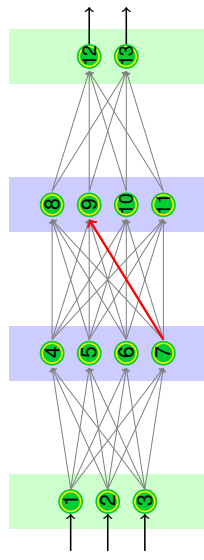
$$o(\vec{x}) = \sigma(\tilde{w} \cdot \vec{x})$$

$$\text{where } \sigma(y) = \frac{1}{1 + e^{-y}}$$

- **Backpropagation** algorithm learns the 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

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

Conventions Over The Network



- x_{ij} i th input to unit j (x_{94} is highlighted)
- w_{ij} weight associated with i th input to unit j
- net_j be $\sum_i w_{ij} x_{ij}$ the weighted sum of input for unit j
- o_j output computed by unit j . Consider it as $\sigma(net_j)$ target output for unit j
- t_j set of units in final layer (**112, 113**) in our case
- $outputs$ units whose immediate input is the output of unit j
- $Downstream(i)$

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

Thank You!

Thank you very much for your attention!