

# Forward Propagation and Backward Propagation

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics (Deep Feedforward Networks)

- Overview
  - 1.Example: Learning XOR
  - 2.Gradient-Based Learning
  - 3.Hidden Units
  - 4.Architecture Design
  - 5.Backpropagation and Other Differentiation Algorithms
  - 6.Historical Notes

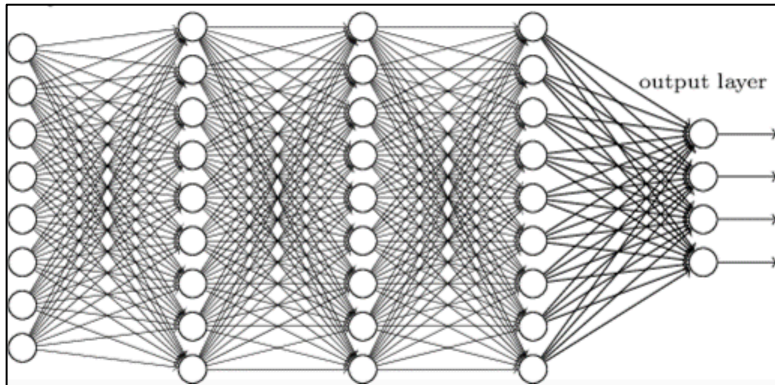
# Topics in Backpropagation

1. Forward and Backward Propagation
2. Computational Graphs
3. Chain Rule of Calculus
4. Recursively applying the chain rule to obtain backprop
5. Backpropagation computation in fully-connected MLP
6. Symbol-to-symbol derivatives
7. General backpropagation
8. Ex: backpropagation for MLP training
9. Complications
10. Differentiation outside the deep learning community
11. Higher-order derivatives

# Forward Propagation

- Producing an output from input
  - When we *use* a Feed-Forward Network to accept an input  $\mathbf{x}$  and produce an output  $\hat{\mathbf{y}}$  information  $\mathbf{x}$  propagates to hidden units at each layer and finally produces  $\hat{\mathbf{y}}$
  - This is called forward propagation
- During training (quality of result is evaluated):
  - forward propagation can continue onward
  - until it produces scalar cost (loss)  $J(\boldsymbol{\theta})$  over  $N$  training samples  $(\mathbf{x}_n, \mathbf{y}_n)$

# Equations for Forward Propagation



## Producing an output:

First layer given by  $\mathbf{h}^{(1)} = g^{(1)}(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)})$

Second layer is  $\mathbf{h}^{(2)} = g^{(2)}(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)})$ ,

.....

Final output is  $\hat{\mathbf{y}} = g^{(d)}(\mathbf{W}^{(d)\top} \mathbf{h}^{(d)} + \mathbf{b}^{(d)})$

During training, cost (loss) over  $N$  exemplars:

$$J(\theta) = J_{MLE} + \lambda \left( \sum_{i,j} (W_{i,j}^{(1)})^2 + \sum_{i,j} (W_{i,j}^{(2)})^2 + \dots \right)$$

$$J_{MLE} = \frac{1}{N} \sum_{n=1}^N \|\hat{\mathbf{y}} - \mathbf{y}_n\|^2$$

# Loss Function

- **Loss function** —  $L$  is used to evaluate to what extent the actual outputs  $y$  are correctly predicted by the model outputs  $z$
- **Cross-entropy loss** — In the context of binary classification in neural networks, the cross-entropy loss  $L(z, y)$  is commonly used:

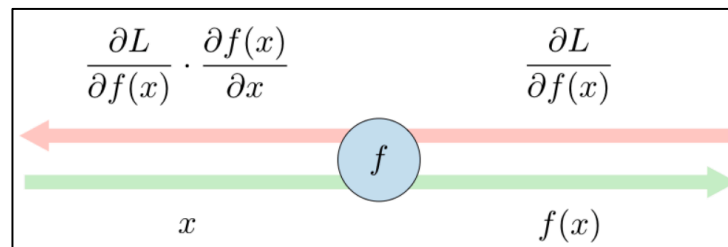
$$L(z, y) = - \left[ y \log(z) + (1 - y) \log(1 - z) \right]$$

# Back-Propagation Algorithm

- Allows information from the cost to flow back to compute gradient (derivative)

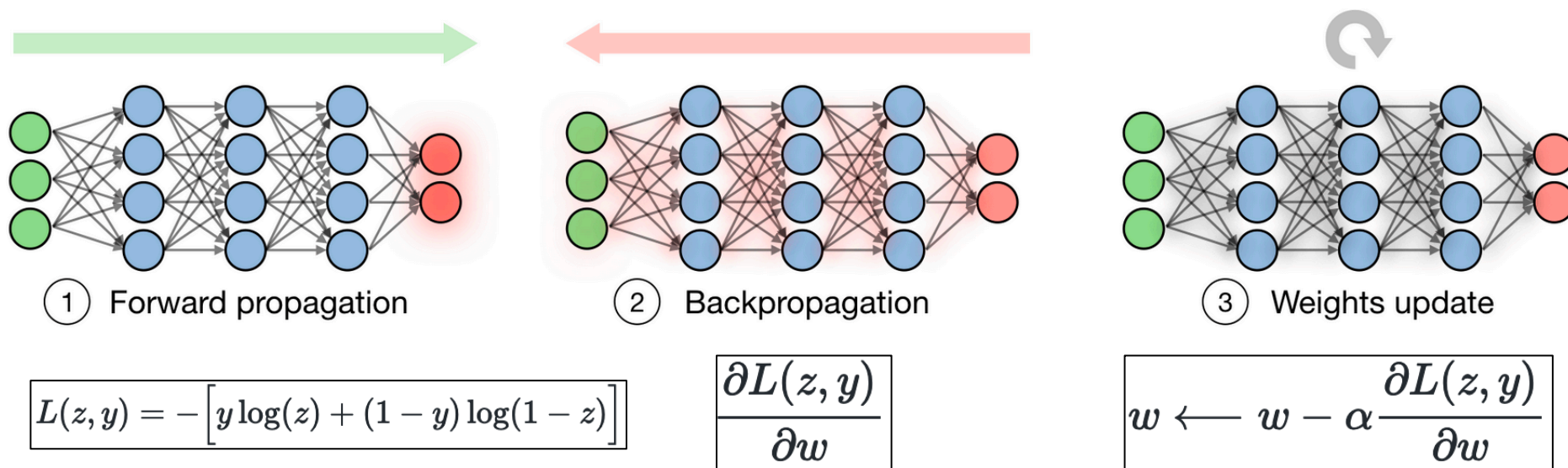
$$\frac{\partial L(z, y)}{\partial w}$$

- Computing analytical expression for the gradient is straightforward
  - But numerically evaluating the gradient is computationally expensive
  - Derivative wrt each weight  $w$  is computed using the chain rule



# Updating Weights

- Step 1: Take a batch of training data and forward propagate to compute the loss
- Step 2: Backpropagate the loss to get the gradient of the loss wrt each weight
- Step 3: Use the gradients to update weights





# Analytical Expression for Gradient

- Gradient is needed to update weights

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla J_n(\theta)$$

- Sum-of-squares criterion over  $N$  samples

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \right\}^2$$

– Expression for gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \sum_{n=1}^N \left\{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \right\} \phi(\mathbf{x}_n)$$

- Another way of saying the same with cost  $J(\theta)$ :

$$J_n(\theta) = \left\| \theta^T \mathbf{x}_n - y_n \right\|^2$$

$$\nabla_{\theta} J_n(\theta) = \theta^T (\theta^T \mathbf{x}_n - y_n) = \theta^T \theta \mathbf{x}_n - \theta^T y_n$$

# Backpropagation is not Learning

- Backpropagation often misunderstood as the whole learning algorithm for multilayer networks
  - It only refers to method of computing gradient
- Another algorithm, e.g., SGD, is used to perform learning using this gradient
  - Learning is updating weights using gradient:
- Backpropagation is also misunderstood to being specific to multilayer neural networks
  - It can be used to compute derivatives for any function (or report that the derivative is undefined)

# Importance of Backpropagation

- Backprop is a technique for computing derivatives quickly
  - It is the key algorithm that makes training deep models computationally tractable
  - For modern neural networks it can make training gradient descent 10 million times faster relative to naive implementation
    - It is the difference between a model that takes a week to train instead of 200,000 years

# Computing gradient for arbitrary function

- Arbitrary function  $f(\mathbf{x}, \mathbf{y})$ 
  - $\mathbf{x}$  : variables for which derivatives  $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$  are desired
  - $\mathbf{y}$  is an additional set of variables that are inputs to the function but whose derivatives are not required
- Gradient required is of cost wrt parameters,  $\nabla_{\theta} J(\theta)$
- Backprop is also useful for other ML tasks
  - Those that need derivatives, as part of learning process or to analyze a learned model
  - To compute Jacobian of a function  $f$  with multiple outputs
- We restrict to case where  $f$  has a single output