

Error Backpropagation

Sargur Srihari

Topics

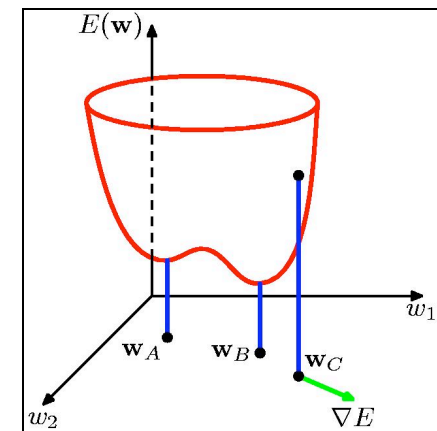
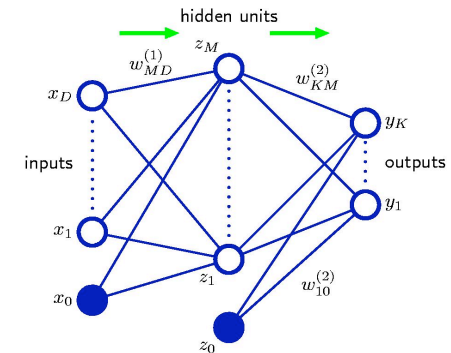
- Neural Network Learning Problem
- Need for computing derivatives of Error function
- Forward propagation of activations
- Backward propagation of errors
- Statement of Backprop algorithm
- Use of backprop in computing the Jacobian matrix

Neural Network Learning Problem

- Goal is to learn the weights \mathbf{w} from a labelled set of training samples
- Learning procedure has two stages
 1. Evaluate derivatives of error function with respect to weights w_1, \dots, w_T
 2. Use derivatives to compute adjustments to weights

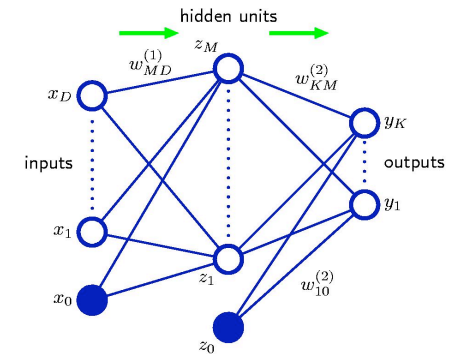
$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})$$

$$\begin{aligned} T &= (D+1)M + (M+1)K \\ &= M(D+K+1) + K \end{aligned}$$



Back-propagation Terminology

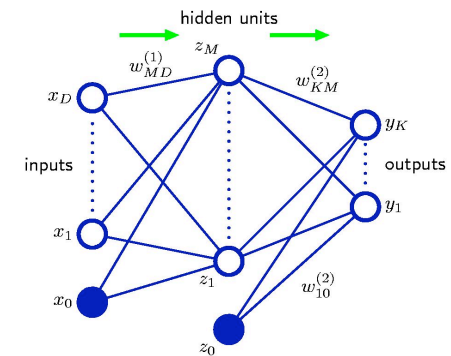
- Goal: Efficient technique for evaluating gradient of an error function $E(\mathbf{w})$ for a feed-forward neural network
- Backpropagation is term used for derivative computation only
- In subsequent stage derivatives are used to make adjustments to weights
- Achieved using a local message passing scheme
 - Information sent forwards and backwards alternately



Overview of Backprop algorithm

- Choose *random* weights for the network
- Feed in an example and obtain a result
- Calculate the error for each node (starting from the last stage and propagating the error backwards)
- Update the weights
- *Repeat* with other examples until the network converges on the target output

- How to divide up the errors needs a little calculus

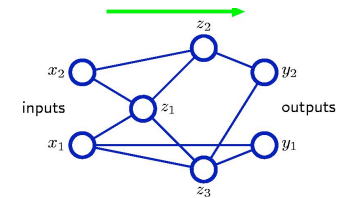


Wide use of Backpropagation

- Can be applied to error function other than sum of squared errors
- Used to evaluate other matrices such as Jacobian and Hessian matrices
- Second stage of weight adjustment using calculated derivatives can be tackled using variety of optimization schemes substantially more powerful than gradient descent

Evaluation of Error Function Derivatives

- Derivation of back-propagation algorithm for
 - Arbitrary feed-forward topology
 - Arbitrary differentiable nonlinear activation function
 - Broad class of error functions
- Error functions of practical interest are sums of errors associated with each training data point



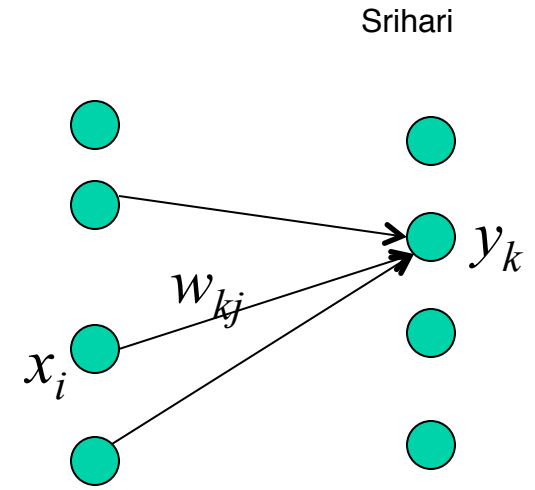
$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

- We consider problem of evaluating
 - For n^{th} term in the error function $\nabla E_n(\mathbf{w})$
 - Derivatives are wrt the weights w_1, \dots, w_T
- Can be used directly for sequential optimization or accumulated over training set (for batch)

A simple Linear Model

- Outputs y_k are linear combinations of input variables x_i

$$y_k = \sum_i w_{ki} x_i$$



- Error function for a particular input \mathbf{x}_n has the form

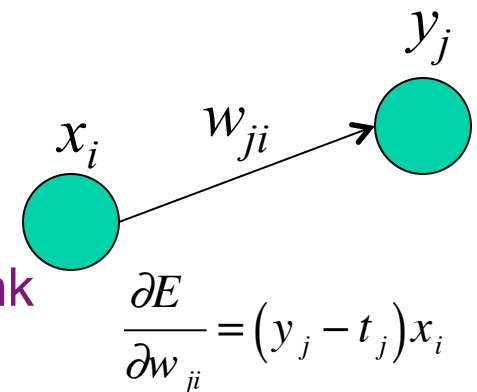
$$E_n = \frac{1}{2} \sum_k (y_{nk} - t_{nk})^2$$

- where $y_{nk} = y_k(\mathbf{x}_n, \mathbf{w})$
- Gradient of Error function wrt a weight w_{ji}

$$\frac{\partial E_n}{\partial w_{ji}} = (y_{nj} - t_{nj}) x_{ni}$$

- a local computation involving product of
 - error signal $y_{nj} - t_{nj}$ associated with output end of link w_{ji}
 - variable x_{ni} associated with input end of link

Subscript n is for a particular input \mathbf{x}_n which is ignored below



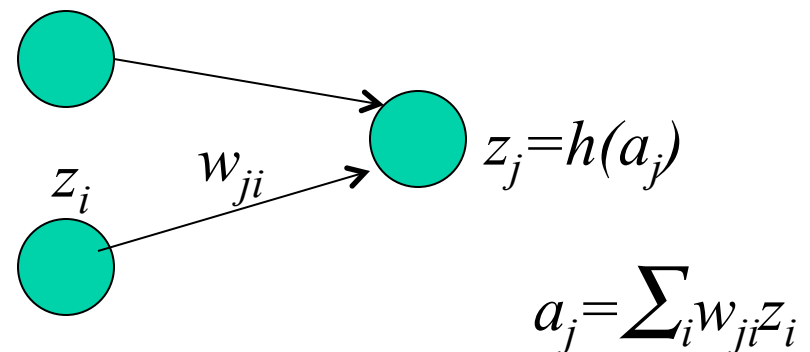
General Feed-Forward Network: Forward Propagation

- Each unit computes weighted sum of its inputs

$$a_j = \sum_i w_{ji} z_i$$

- z_i is activation of a unit (or input) that sends a connection to unit j and w_{ji} is the weight associated with the connection
- Transformed by nonlinear activation function

$$z_j = h(a_j)$$



Evaluation of Derivative E_n wrt a weight w_{ji}

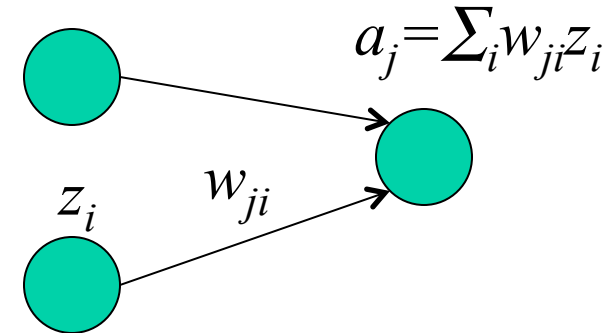
- By chain rule for partial derivatives

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}}$$

Define $\delta_j \equiv \frac{\partial E_n}{\partial a_j}$

$$a_j = \sum_i w_{ji} z_i$$

we have $\frac{\partial a_j}{\partial w_{ji}} = z_i$



- Substituting

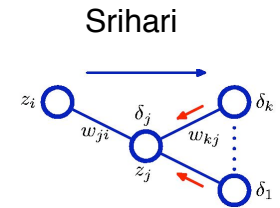
$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$$

- Thus required derivative $\frac{\partial E_n}{\partial w_{ji}}$ is obtained by

- Multiplying value of δ for the unit at output end of weight by value of z for unit at input end of weight

- Need to figure out how to calculate $\delta_j \equiv \frac{\partial E_n}{\partial a_j}$

Calculation of Error for hidden unit δ_j



- For output unit $\delta_k = y_k - t_k$

Since $E = \frac{1}{2} \sum_k (y_k - t_k)^2$ and $y_k = a_k = \sum w_{ki} z_i$ $\delta_k \equiv \frac{\partial E}{\partial a_k}$

- For hidden unit j

- By chain rule

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \underbrace{\frac{\partial E_n}{\partial a_k}}_{\delta_k} \underbrace{\frac{\partial a_k}{\partial a_j}}_{w_{kj}}$$

We are summing partial derivatives over several variables a_k

- Substituting

- We get the backpropagation formula for error derivatives at stage j

$$\delta_k \equiv \frac{\partial E_n}{\partial a_k}$$

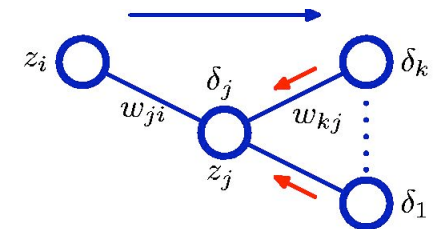
$$a_k = \sum_i w_{ki} z_i = \sum_i w_{ki} h(a_i)$$

$$\frac{\partial a_k}{\partial a_j} = \sum_k w_{kj} h'(a_j)$$

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

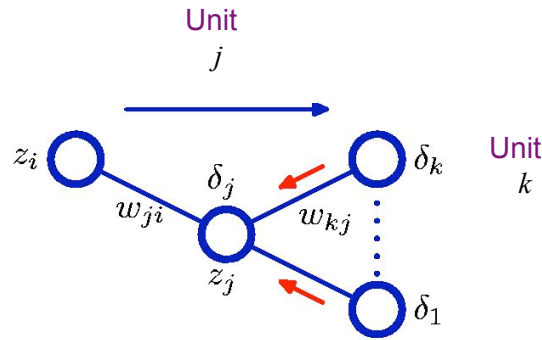
Input to activation from earlier units

error derivative at later unit k



Blue arrow for forward propagation
Red arrows indicate direction of information flow during error backpropagation

Error Backpropagation Algorithm



- Backpropagation Formula

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

- Value of δ for a particular hidden unit can be obtained by propagating the δ 's backward from units higher-up in the network

1. Apply input vector x_n to network and forward propagate through network using

$$a_j = \sum_i w_{ji} z_i \quad \text{and} \quad z_j = h(a_j)$$

2. Evaluate δ_k for all output units using $\delta_k = y_k - t_k$

3. Backpropagate the δ 's using

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

to obtain δ_j for each hidden unit

4. Use $\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$

to evaluate required derivatives

A Simple Example

- Two-layer network
- Sum-of-squared error
- Output units: *linear activation functions*

$$y_k = a_k$$

- Hidden units have *logistic sigmoid* activation function

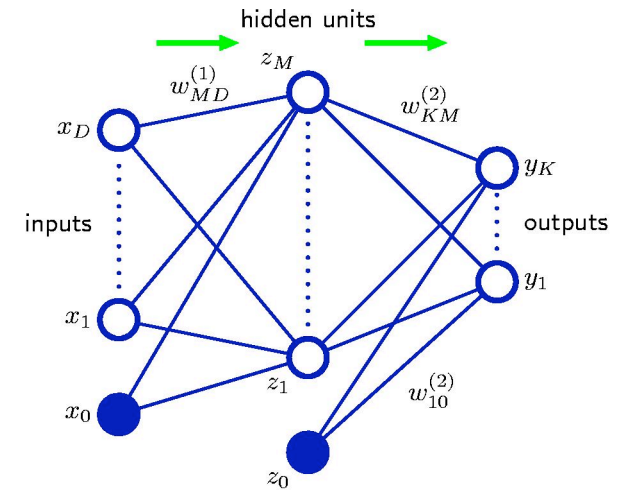
$$h(a) = \tanh(a)$$

where

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

simple form for derivative

$$h'(a) = 1 - h(a)^2$$



Standard Sum of Squared Error

$$E_n = \frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2$$

y_k : activation of output unit k
 t_k : corresponding target
 for input x_k

Simple Example: Forward and Backward Prop

For each input in training set:

- Forward Propagation

$$\left\{ \begin{array}{l} a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i \\ z_j = \tanh(a_j) \\ y_k = \sum_{j=0}^M w_{kj}^{(2)} z_j \end{array} \right.$$

- Output differences

$$\delta_k = y_k - t_k$$

- Backward Propagation (δ s for hidden units)

$$\delta_j = (1 - z_j^2) \sum_{k=1}^K w_{kj} \delta_k$$

- Derivatives wrt first layer and second layer weights

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \delta_j x_i \quad \frac{\partial E_n}{\partial w_{kj}^{(2)}} = \delta_k z_j$$

- Update the weights using $w^{(\tau+1)} = w^{(\tau)} - \eta \nabla E_n(w^{(\tau)})$

Efficiency of Backpropagation

- Computational Efficiency is main aspect of back-prop
- Number of operations to compute derivatives of error function scales with total number W of weights and biases
- Single evaluation of error function for a single input requires $O(W)$ operations (for large W)
- This is in contrast to $O(W^2)$ for numerical differentiation
 - As seen next

Another Approach: Numerical Differentiation

- Compute derivatives using method of finite differences
 - Perturb each weight in turn and approximate derivatives by

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \varepsilon) - E_n(w_{ji})}{\varepsilon} + O(\varepsilon) \text{ where } \varepsilon \ll 1$$

- Accuracy improved by make ε smaller until round-off problems arise
- Accuracy can be improved by using central differences

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \varepsilon) - E_n(w_{ji} - \varepsilon)}{2\varepsilon} + O(\varepsilon^2)$$

- This is $O(W^2)$
- Useful to check if software for backprop has been correctly implemented (for some test cases)

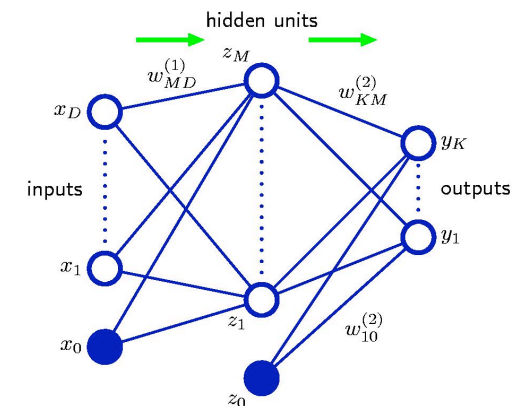
Summary of Backpropagation

- Derivatives of error function wrt weights are obtained by propagating errors backward
- It is more efficient than numerical differentiation
- It can also be used for other computations
 - As seen next for Jacobian

The Jacobian Matrix

- For a vector valued output $\mathbf{y} = \{y_1, \dots, y_m\}$ with vector input $\mathbf{x} = \{x_1, \dots, x_n\}$,
- Jacobian matrix organizes all the partial derivatives into an $m \times n$ matrix

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad J_{ki} = \frac{\partial y_k}{\partial x_i}$$



For a neural network we have a $D+1$ by K matrix

Determinant of Jacobian Matrix is referred to simply as the Jacobian

Jacobian Matrix Evaluation

- In backprop, derivatives of error function wrt weights obtained by propagating errors backward
 - We calculate $\frac{\partial E_n}{\partial w_{ji}} = (y_{nj} - t_{nj})x_{ni}$
- Method can also be used to calculate other derivatives
- Evaluation of Jacobian matrix
 - Elements of matrix are derivatives of network outputs wrt inputs

$$J_{ki} = \frac{\partial y_k}{\partial x_i}$$

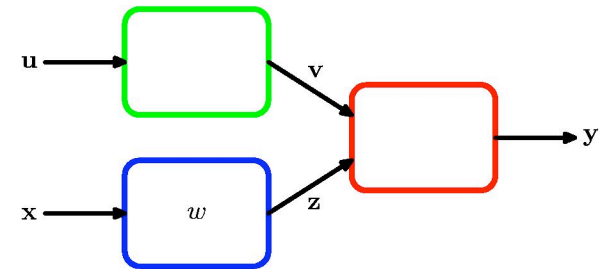
Note that for backprop we are taking the output to be $y_k - t_k$ which is the error

- Each derivative evaluated with other inputs fixed

Use of Jacobian Matrix

- Jacobian plays useful role in systems built from several modules
- We wish to minimize error E wrt parameter w

$$\frac{\partial E}{\partial w} = \sum_{k,j} \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_j} \frac{\partial z_j}{\partial w}$$



- Jacobian matrix for red module appears in the middle term
- Jacobian matrix provides measure of sensitivity of outputs to changes in input

Jacobian Matrix Computation

- Apply input vector corresponding to point in input space where the Jacobian matrix is to be found
- Forward propagate to obtain activations of the hidden and output units in the network
- For each row k of the Jacobian matrix, corresponding to output unit k :
 - Backpropagate for all the hidden units in the network
 - Finally backpropagate to the inputs
- Implementation of such an algorithm can be checked using numerical differentiation

Summary

- Neural network learning needs learning of weights from samples
- Involves two steps: determining derivative of output of a unit wrt each input
- Backpropagation is a general term for computing derivatives
 - Evaluate δ_k for all output units using $\delta_k = y_k - t_k$
 - Backpropagate the δ_k 's to obtain δ_j for each hidden unit
 - Product of δ 's with activations at the unit provide the derivatives for that weight
- Useful to compute a Jacobian matrix with several inputs and outputs