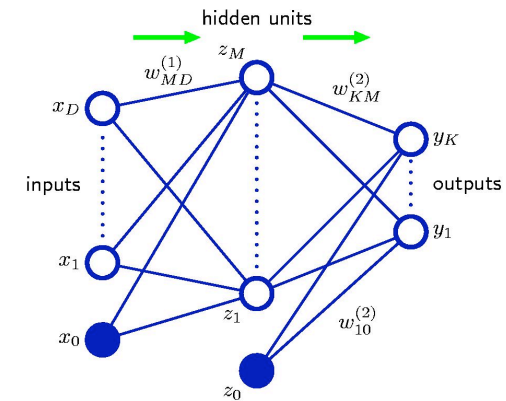


# The Hessian Matrix

Sargur Srihari

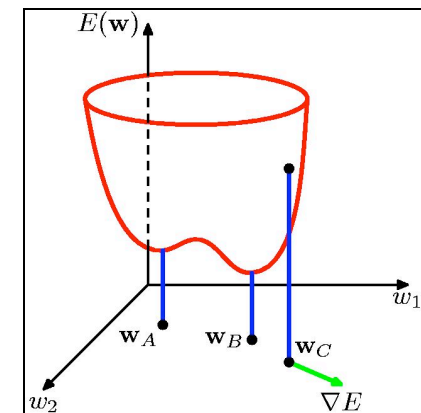
## Hessian of Neural Network Error Function

- Backpropagation can be used to obtain first derivatives of error function wrt weights in network
- It can also be used to derive second derivatives



$$\frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}}$$

- If all weights and bias parameters are elements  $w_i$  of single vector  $\mathbf{w}$  then the second derivatives form the elements  $H_{ij}$  of Hessian matrix  $\mathbf{H}$  where  $i, j \in \{1, \dots, W\}$



# Hessian Matrix Definition

If  $f$  is a real-valued function  $f(x_1, \dots, x_n)$

And if all real-valued derivatives exist then the Hessian matrix of  $f$  is the matrix

$H(f)_{ij}(\mathbf{x}) = D_i D_j f(\mathbf{x})$  where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $D_i$  is the differential operator wrt the  $i^{\text{th}}$  variable

$$H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Hessian matrices are used in large-scale optimization problems within Newton-type methods because they are the coefficient of the quadratic term of a local Taylor expansion of a function.

That is,

$$y = f(\mathbf{x} + \Delta \mathbf{x}) \approx f(\mathbf{x}) + J(\mathbf{x})\Delta \mathbf{x} + \frac{1}{2}\Delta \mathbf{x}^T H(\mathbf{x})\Delta \mathbf{x}$$

where  $J$  is the Jacobian matrix, which is a vector (the gradient) for scalar-valued functions.

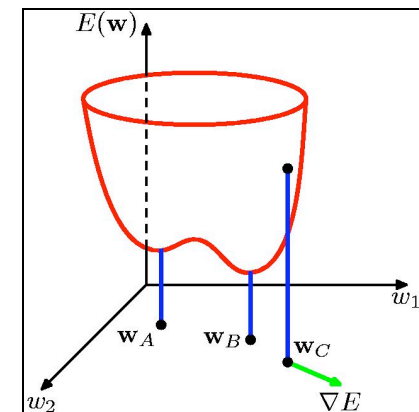
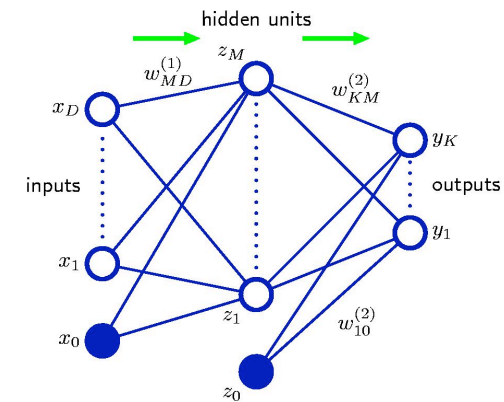
For a vector-valued function  $\mathbf{y} = (y_1, \dots, y_n)$

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

In backprop we deal with scalar error function  $E$

# Role of Hessian in Neural Computing

1. Several nonlinear optimization algorithms for neural networks  
are based on second order derivatives of error surface
2. Basis for fast procedure for retraining with small change of training data
3. Identifying least significant weights  
for network pruning requires inverse of Hessian
4. Bayesian neural network  
Central role in Laplace approximation



# Evaluating the Hessian Matrix

- Full Hessian matrix can be difficult to compute in practice
  - quasi-Newton algorithms have been developed that use approximations to the Hessian
- Various approximation techniques have been used to evaluate the Hessian for a neural network
  - calculated exactly using an extension of backpropagation
- Important consideration is efficiency
  - With  $W$  parameters (weights and biases) matrix has dimension  $W \times W$
  - Efficient methods have  $O(W^2)$

## Methods for evaluating the Hessian Matrix

- Diagonal Approximation
- Outer Product Approximation
- Inverse Hessian
- Finite Differences
- Exact Evaluation using Backpropagation
- Fast multiplication by the Hessian

# Diagonal Approximation

- In many case inverse of Hessian is needed
- If Hessian approximation is diagonal, its inverse is trivially computed
- Complexity is  $O(W)$  rather than  $O(W^2)$  for full Hessian

# Outer product approximation

- Neural networks commonly use sum-of-squared errors function

$$E = \frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2$$

- Can write Hessian matrix in the form

$$H \approx \sum_{n=1}^n b_n b_n^T$$

- Where  $b_n = \nabla y_n = \nabla a_n$
- Elements can be found in  $O(W^2)$  steps

# Inverse Hessian

- Use outer product approximation to obtain computationally efficient procedure for approximating inverse of Hessian

# Finite Differences

- Using backprop, complexity is reduced from  $O(W^3)$  to  $O(W^2)$

# Exact Evaluation of the Hessian

- Using an extension of backprop
- Complexity is  $O(W^2)$

# Fast Multiplication by the Hessian

- Application of the Hessian involve multiplication by the Hessian
- The vector  $\mathbf{v}^T \mathbf{H}$  has only  $W$  elements
- Instead of computing  $\mathbf{H}$  as an intermediate step, find efficient method to compute product