

The Hessian Matrix

Sargur Srihari

Definitions of Gradient and Hessian

- First derivative of a scalar function $E(\mathbf{w})$ with respect to a vector $\mathbf{w}=[w_1, w_2]^T$ is a vector called the *Gradient* of $E(\mathbf{w})$

$$\nabla E(\mathbf{w}) = \frac{d}{d\mathbf{w}} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \end{bmatrix}$$

If there are M elements in the vector then Gradient is a $M \times 1$ vector

- Second derivative of $E(\mathbf{w})$ is a matrix called the *Hessian* of $E(\mathbf{w})$

$$H = \nabla \nabla E(\mathbf{w}) = \frac{d^2}{d\mathbf{w}^2} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2} & \frac{\partial^2 E}{\partial w_1 \partial w_2} \\ \frac{\partial^2 E}{\partial w_2 \partial w_1} & \frac{\partial^2 E}{\partial w_2^2} \end{bmatrix}$$

Hessian is a matrix with M^2 elements

- Jacobian* is a matrix consisting of first derivatives wrt a vector

Need for Gradient and Hessian in ML

Training samples: $n=1, \dots, N$

Inputs: $M \times 1$ vectors

$$\phi(\mathbf{x}_n) = \begin{pmatrix} \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \dots & \phi_{M-1}(\mathbf{x}_n) \end{pmatrix}^T$$

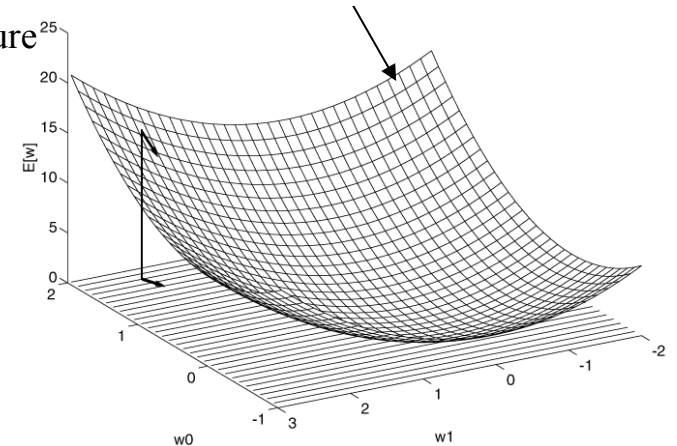
Outputs

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$\Phi_0(\mathbf{x})=1$, dummy feature

Error surface for $M=2$
a paraboloid with
a single global minimum



For Logistic Regression (cross-entropy error):

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\} \text{ where } y_n = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$$

For Stochastic Gradient Descent we need $\nabla E_n(\mathbf{w})$ where $E(\mathbf{w}) = \sum_n E_n(\mathbf{w})$

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

For Newton-Raphson update we need both $\nabla E(\mathbf{w})$ and $H = \nabla \nabla E(\mathbf{w})$

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

Gradient and Hessian for Two-class Logistic Regression

- Cross-Entropy Error

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\} \text{ where } y_n = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$$

- Gradient of E

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi(\mathbf{x}_n) = \Phi^T (\mathbf{y} - \mathbf{t})$$

$$\mathbf{y} = (y_1, \dots, y_N)^T$$

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

- Hessian of E

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi(\mathbf{x}_n) \phi^T(\mathbf{x}_n) = \Phi^T \mathbf{R} \Phi$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

\mathbf{R} is $N \times N$ diagonal matrix with elements

$$R_{nn} = y_n (1 - y_n) = \mathbf{w}^T \phi(\mathbf{x}_n) (1 - \mathbf{w}^T \phi(\mathbf{x}_n))$$

$$\phi(\mathbf{x}_n) = \left(\phi_0(\mathbf{x}_n) \quad \phi_1(\mathbf{x}_n) \quad \dots \quad \phi_{M-1}(\mathbf{x}_n) \right)^T$$

Hessian is not constant and depends on \mathbf{w} through \mathbf{R}

Since \mathbf{H} is positive-definite (i.e., for arbitrary \mathbf{u} , $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$)

error function is a concave function of \mathbf{w} and so has a unique minimum

Gradient and Hessian for Multi-class Logistic Regression

- Cross-Entropy Error

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

$$T = \begin{bmatrix} t_{11} & \cdot & \cdot & t_{1K} \\ \cdot & & & \\ \cdot & & & \\ t_{N1} & & & t_{NK} \end{bmatrix}$$

- Gradient of E

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N (y_{nj} - t_{nj}) \phi(\mathbf{x}_n)$$

$$y_{nk} = y_k(\phi(\mathbf{x}_n)) = \frac{\exp(\mathbf{w}_k^T \phi(\mathbf{x}_n))}{\sum_j \exp(\mathbf{w}_j^T \phi(\mathbf{x}_n))}$$

- Hessian of E

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi(\mathbf{x}_n) \phi^T(\mathbf{x}_n)$$

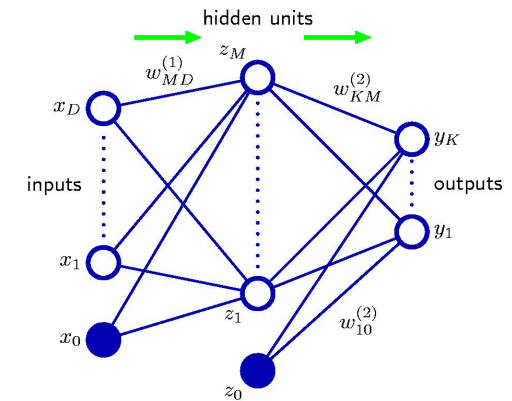
$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$$\phi(\mathbf{x}_n) = \left(\phi_0(\mathbf{x}_n) \quad \phi_1(\mathbf{x}_n) \quad \dots \quad \phi_{M-1}(\mathbf{x}_n) \right)^T$$

Each element of the Hessian needs M multiplications and additions
 Since there are M^2 elements in the matrix the computation is $O(M^3)$

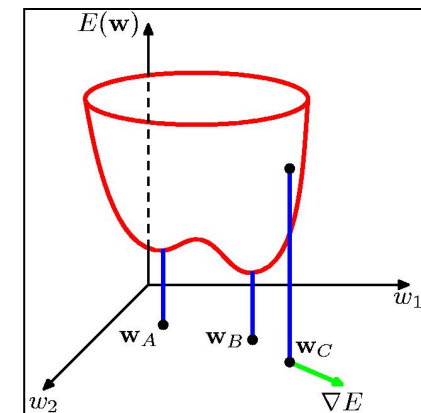
Hessian of Neural Network Error Function

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



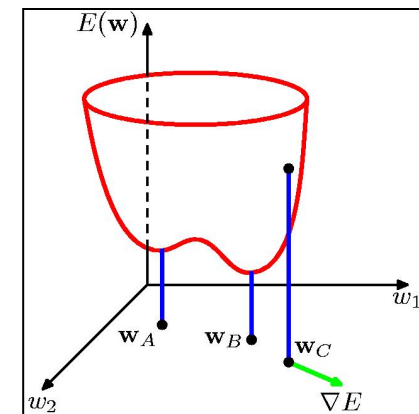
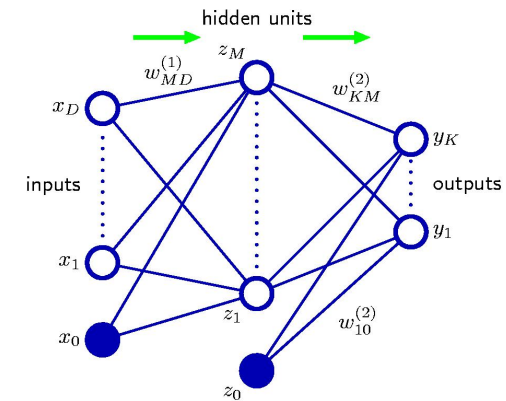
$$\frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}}$$

- 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 $i, j \in \{1, \dots, W\}$



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