

## REVISING THE TRADITIONAL BACKPROPAGATION WITH THE METHOD OF VARIABLE METRIC(QUASI - NEWTON) AND APPROXIMATING A STEP SIZE

Sang Woong CHOE <sup>a</sup> , Jin Choon LEE <sup>b</sup>

a. Dept. of Management Information System, Pohang Junior College  
 Hunghae-Up, Pohang, Kyungpook, 791-940, Korea  
 Tel: +82-562-45-1198

b. Dept. of Industrial Engineering, Kyungil University  
 Hayang, Kyungsan, Kyungpook, 712-701, Korea  
 Tel: +82-53-850-7266, Fax: +82-53-850-7275, E-mail: jinlee@bear.kyungil.ac.kr

### Abstract

In this paper, we propose another paradigm(QNBP) to be capable of overcoming limitations of the traditional backpropagation(SDBP). QNBP is based on the method of Quasi-Newton(variable metric) with the normalized direction vectors and computes step size through the linear search. Simulation results showed that QNBP was definitely superior to both the stochastic SDBP and the deterministic SDBP in terms of accuracy and rate of convergence and might surmount the problem of local minima. And there was no difference between DFP+SR1 and BFGS+SR1 combined algorithms in QNBP.

**Keywords :** Quasi-Newton(variable metric) method, combined algorithm, DFP+SR1, BFGS+SR1, normalized direction vector

### 1. Introduction

The traditional backpropagation method based on the steepest descent method(here-after referred to as the SDBP) is not good learning algorithm in terms of accuracy and rate of convergence. Especially, the most serious problem in SDBP training is the existence of local minima, where the error at the network outputs may still be unacceptably high. It is undeniable that most researchers have been inclined to blindly accept this method for many subsequent years.

In this paper, we propose another paradigm to be capable of overcoming limitations of SDBP. The proposed backpropagation (here-after referred to as the QNBP) is based on the combined Quasi-Newton method(or variable metric method) with the approximate to step sizes. The term "combined" implies that two rank-two updates in Quasi-Newton methods can be combined with the symmetric rank-one update.

To use the Newton-Raphson method we must know the Hessian matrix of the function and we must invert it. In many cases including backpropagation model, the Hessian matrices are not available, are very expensive to compute. The variable metric method, originally developed by Davidon[1], avoids these difficulties.

At current point  $x_k \in R^n$ , arbitrary function  $f(x)$  can be locally approximated by the quadratic form of Eq. (1) using Taylor series

$$f(x_k) \approx b + a^T x_k + \frac{1}{2} x_k^T H(x_k) x_k, H(x_k) > 0 \quad (1)$$

where  $H(x_k)$  is the Hessian matrix of  $f(x_k)$ .

$f(x_k)$  has the minimum point  $x_{k+1} = -H(x_k)^{-1}a$ . From the

gradient vector  $g_k$ ,  $x_{k+1}$  can be computed in one step as

$$\begin{aligned} x_{k+1} &= -H(x_k)^{-1}a = -H(x_k)^{-1}(g_k - H(x_k)x_k) \\ &= -H(x_k)^{-1}g_k + x_k \end{aligned} \quad (2)$$

if  $H(x_k)^{-1}$  is available.

The optimal direction of descent is along  $-H(x_k)^{-1}g_k$  and not along  $-g_k$ . And  $H(x_k)^{-1}g_k$  is the gradient of  $f(x_k)$  when the metric is given by  $\sqrt{x_k^T H(x_k) x_k}$ . In most cases,  $H(x_k)$  is not known. Therefore, this method estimates  $H(x_k)^{-1}$  successively by constructing a sequence of matrices such that they are the inverse of  $H(x_k)$ . Hence, the name of the variable metric method or Quasi-Newton method.

The "Quasi" in Quasi-Newton method is because we don't use the actual Hessian matrix, but instead use current approximation of it. This is often better than using the true Hessian. In short, the basic idea of the Quasi-Newton methods is to build up, iteratively, a good estimate to the inverse Hessian matrix.

Quasi-Newton methods come in three main flavors; two rank-two algorithms such as Davidon - Fletcher - Powell(DFP) algorithm[2], Broyden - Fletcher - Goldfarb - Shanno(BFGS) algorithm[3,4,5,6] and symmetric rank-one algorithm(SR1)[7]. This three versions differ mainly in choice of an estimate to the inverse Hessian matrix. The main drawback of SR1 update is that an estimate to inverse Hessian at next point may not be positive definite even when an estimate to inverse Hessian at current point is positive definite.

However, the main advantage of SR1 update is that it requires less amount of storage space and less amount of arithmetic operations in

comparison to two rank-two updates such as DFP and BFGS updates. Hence, if the fatal flaw of SR1 update could be overcome, this update is the most adequate to the backpropagation model.

In this paper, two combined updates, DFP+SR1 and BFGS+SR1, are applied to the backpropagation model. So, we intend to call this another paradigm QNBP.

## 2. QNBP : Algorithm and step size

First, to describe QNBP algorithm, we introduce the following notation.

- [1]  $r$  -- \* layer : input layer(0), hidden layer(1~(L-1)), and outputlayer(L).
- [2]  $p, k$  -- \* the  $p$ \_th training data, the  $k$ \_th iteration. :  $1 \leq p \leq P, k \geq 0$
- [3]  $N_r$  -- \* number of the  $r$ -layer nodes except bias term.  
 $n_r$  -- \* the  $n_r$ \_th  $r$ -layer node. :  $0 \leq n_r \leq N_r, (r \neq L)$   
and bias term, if  $n_r = N_r$ .
- [4]  $x_{pn_p}$  -- \* value of the  $p$ \_th training data to the  $n_0$ \_th input layer node.  
 $y_{pn_L}$  -- \* value of the  $p$ \_th training data to the  $n_L$ \_th output layer node.
- [5]  $w_{n_r, n_{r-1}}^r, d_{n_r, n_{r-1}}^r$  -- \* weight and direction on the connection from the  $n_{r-1}$ \_th ( $r-1$ )-layer node to the  $n_r$ \_th  $r$ -layer node. :  $n_r \neq N_r$
- [6]  $net_{pn_r}^r, f_{n_r}^r(net_{pn_r}^r)$  -- \* net value and output value of the  $n_r$ \_th  $r$ -layer node for the  $p$ \_th training data.
- [7]  $f_{n_0}^0(net_{pn_0}^0) = net_{pn_0}^0 = x_{pn_0}$  -- \*  $n_0 \neq N_0$   
 $f_{n_r}^r(net_{pn_r}^r) = f_{n_r}^r(0) = 1$  -- \*  $n_r = N_r$   
 $f_{n_r}^r(net_{pn_r}^r) = f_{n_r}^r(\sum_{n_{r-1}} w_{n_r, n_{r-1}}^r f_{n_{r-1}}^{r-1}(net_{pn_{r-1}}^{r-1}))$   
 $= f_{n_r}^r(\sum_{n_{r-1}=N_{r-1}} (w_{n_r, n_{r-1}}^r f_{n_{r-1}}^{r-1}(net_{pn_{r-1}}^{r-1})) + w_{n_r, N_{r-1}}^r)$   
-- \*  $r \neq 0, n_r \neq N_r$
- [8]  $g_{n_r, n_{r-1}}^r = \frac{\partial E(\cdot)}{\partial w_{n_r, n_{r-1}}^r} = -\sum_p \delta_{pn_r}^r f_{n_{r-1}}^{r-1}(net_{pn_{r-1}}^{r-1})$   
 $\delta_{pn_r}^r = \frac{df_{n_r}^r(net_{pn_r}^r)}{dnet_{pn_r}^r} \sum_{n_{r-1}} \delta_{pn_{r-1}}^{r-1} w_{n_r, n_{r-1}}^{r-1}, r \neq 0, L$  and  
 $\delta_{pn_L}^L = \frac{df_{n_L}^L(net_{pn_L}^L)}{dnet_{pn_L}^L} (y_{pn_L} - f_{n_L}^L(net_{pn_L}^L))$
- [9]  $W^r = [w_{n_r, n_{r-1}}^r]_{N_r \times (N_{r-1} + 1)}, D^r = [d_{n_r, n_{r-1}}^r]_{N_r \times (N_{r-1} + 1)}$   
 $G^r = [g_{n_r, n_{r-1}}^r]_{N_r \times (N_{r-1} + 1)}$   
-- \*  $r$ -layer weight, direction and gradient matrix  
:  $r \neq 0, n_r \neq N_r$
- [10]  $Vec(A) = [a_{ij}]_{I \times J} = [b_h]_{I \times 1}$  -- \*  $h = jI + i$ , where  $j$  is a maximum integer less than or equal to  $\frac{h}{I}$  and  $0 \leq h \leq (IJ-1), 0 \leq i \leq (I-1), 0 \leq j \leq (J-1)$ .

$$[11] E(\theta) = E(\text{Vec}(W^1), \text{Vec}(W^2), \dots, \text{Vec}(W^L))$$

$$= \frac{1}{2} \sum_p \sum_{n_L} [y_{pn_L} - f_{n_L}^L(\sum_{n_{L-1}} w_{n_L, n_{L-1}}^L f_{n_{L-1}}^{L-1}(\sum_{n_{L-2}} w_{n_{L-1}, n_{L-2}}^{L-1} f_{n_{L-2}}^{L-2}(\dots \dots \dots w_{n_1, n_0}^1 f_{n_0}^0(\sum_{n_1} w_{n_1, n_0}^1 f_{n_0}^0(\text{net}_{pn_0}^0)))))]^2$$

-- \* objective(error) function of the network. :  
 $w_{n_r, n_{r-1}}^r = 0 (r \neq 0, L)$

$$[12] \nabla E(\theta) = \rho = \left[ \frac{\partial E(\theta)}{\partial \text{Vec}(W^r)} \right]_{L \times 1}, \quad \frac{\partial E(\theta)}{\partial \text{Vec}(W^r)} = \text{Vec}(G^r)$$

-- \* gradient vector of  $E(\theta)$ .

$$[13] H(\theta) = [H_{ab}]_{L \times L}, \quad H_{ab} = \frac{\partial^2 E(\theta)}{\partial \text{Vec}(W^b) \partial \text{Vec}(W^a)}$$

$$= \left[ \frac{\partial^2 E(\theta)}{\partial w_{n_a, n_{a-1}}^a \partial w_{n_b, n_{b-1}}^b} \right]_{p \times q}$$

\* -- \* positive definite Hessian matrix. :

$$p = N_a \times (N_{a-1} + 1), \quad q = N_b \times (N_{b-1} + 1)$$

$$[14] \|x\| = (x^T x)^{\frac{1}{2}}$$

$$[15] \theta^r = [\text{Vec}(W^1), \dots, \text{Vec}(W^r), \text{Vec}(W^{r+1}), \dots, \text{Vec}(W^L)]^T$$

$$\lambda^r = [\text{Vec}(D^1), \dots, \text{Vec}(D^r), \text{Vec}(D^{r+1}), \dots, \text{Vec}(D^L)]^T$$

$$\rho^r = [\text{Vec}(G^1), \dots, \text{Vec}(G^r), \text{Vec}(G^{r+1}), \dots, \text{Vec}(G^L)]^T$$

$$\sum_r \theta^r = \theta, \quad \sum_r \lambda^r = \lambda, \quad \sum_r \rho^r = \rho, \quad r \neq 0$$

$$[16] \langle \theta^r \rangle^T = [\text{Vec}^T(W^1), \dots, \text{Vec}^T(W^r), \text{Vec}^T(W^{r+1}), \dots, \text{Vec}^T(W^L)]$$

$$\langle \lambda^r \rangle^T = [\text{Vec}^T(D^1), \dots, \text{Vec}^T(D^r), \text{Vec}^T(D^{r+1}), \dots, \text{Vec}^T(D^L)]$$

$$\langle \rho^r \rangle^T = [\text{Vec}^T(G^1), \dots, \text{Vec}^T(G^r), \text{Vec}^T(G^{r+1}), \dots, \text{Vec}^T(G^L)]$$

$$\sum_r \langle \theta^r \rangle^T = \theta^T, \quad \sum_r \langle \lambda^r \rangle^T = \lambda^T, \quad \sum_r \langle \rho^r \rangle^T = \rho^T, \quad r \neq 0$$

### 2.1 Algorithm

Let  $A_k$  be an estimate of  $H(x_k)^{-1}$  at current point  $\theta_k$ . In QNBP, the next point  $\theta_{k+1}$  is generated from  $\theta_k$  by the linear search in the direction of  $\lambda_k$  and is expressed as

$$\theta_{k+1} = \theta_k + \tau_k \lambda_k, \quad k \geq 0$$

$$\text{where } \lambda_k = \frac{1}{\|A_k \rho_k\|} (-A_k \rho_k) \quad (3)$$

$\tau_k$  is a step size and we select it such that

$$\text{Min}_\tau E(\theta_{k+1})$$

$$\cong \text{Min}_\tau \left\{ E(\theta_k) + \langle \rho_k \rangle^T \lambda_k \tau_k + \frac{1}{2} \langle \lambda_k \rangle^T H(\theta_k) \lambda_k (\tau_k)^2 \right\} \quad (4)$$

From Eq. (4),

$$\tau_k = \frac{-\langle \rho_k \rangle^T \lambda_k}{\langle \lambda_k \rangle^T H(\theta_k) \lambda_k} > 0 \quad (5)$$

Eq. (5) can be written as

$$\tau_k \cong \frac{2(E(\theta_{k+1}) - E(\theta_k))}{\langle \rho_{k+1} \rangle^T \lambda_k + \langle \rho_k \rangle^T \lambda_k} \quad (6)$$

using Taylor series and  $\frac{1}{\tau_{k-1}}(\rho_k - \rho_{k-1}) \cong H(\theta_{k-1})\lambda_{k-1}$ .

One important feature of QNBP is the choice of matrices  $A_k$  to be positive definite and satisfying the following Quasi-Newton equation

$$A_{k+1}p_k = q_k \quad (7)$$

where  $p_k = \rho_{k+1} - \rho_k$  and  $q_k = \theta_{k+1} - \theta_k$ .

QNBP with DFP, BFGS and SR1 updates generates the sequence of  $A_k (k \geq 1)$  respectively as

$$A_{k+1} = A_k + \frac{q_k q_k^T}{q_k^T p_k} - \frac{A_k p_k p_k^T A_k}{p_k^T A_k p_k} \quad (8)$$

$$A_{k+1} = A_k + \frac{1}{q_k^T p_k} \left( \left( 1 + \frac{p_k^T A_k p_k}{q_k^T p_k} \right) q_k q_k^T - q_k p_k^T A_k - A_k q_k p_k^T \right) \quad (9)$$

$$A_{k+1} = A_k + \frac{(q_k - A_k p_k)(q_k - A_k p_k)^T}{(q_k - A_k p_k)^T p_k} \quad (10)$$

First combined algorithm, DFP+SR1, computes  $A_{k+1}$  as follows

$$\begin{aligned} A_{k+1} &= \text{Eq. (10)} && \text{if } (q_k - A_k p_k)^T p_k > 0 \\ &= \text{Eq. (8)} && \text{if } (q_k - A_k p_k)^T p_k < 0 \text{ and } q_k^T p_k > 0 \\ &= sI && \text{otherwise} \end{aligned} \quad (11)$$

where  $I$  is identity matrix and  $s$  is a given positive constant.

Second combined algorithm, BFGS+SR1, computes  $A_{k+1}$  as follows

$$\begin{aligned} A_{k+1} &= \text{Eq. (10)} && \text{if } (q_k - A_k p_k)^T p_k > 0 \\ &= \text{Eq. (9)} && \text{if } (q_k - A_k p_k)^T p_k < 0 \text{ and } q_k^T p_k > 0 \\ &= sI && \text{otherwise} \end{aligned} \quad (12)$$

And  $A_0 = I$

Note that Eq. (11) and (12) check whether the sequence of  $A_k (k \geq 1)$  are positive definite. Thus we are led to consider an algorithm, that is, QNBP which have two versions, DFP+SR1 and BFGS+SR1.

## 2.2 Step size

Eq. (6) may not be a directly available expression since we must estimate  $\rho_{k+1}$  and  $E(\theta_{k+1})$ . For this reason, we will grant two desirable properties to  $\tau_k$  with the purpose of a successful training.

We can analytically derive the unique  $\tau_k$  which satisfies two properties and then use this as the good approximate to a true  $\tau_k$ . These two desirable properties are as follows.

[Property 1]  $\theta_{k+1}$  lies in a given direction  $\lambda_k$  from  $\theta_k$ .

[Property 2]  $E(\theta_k)$  decrease log-linearly.

In short, Eq. (6) satisfies the above-mentioned properties, it can be written as

$$\tau_k \cong \frac{2(E_1 - E_0) \left( \frac{E_1}{E_0} \right)^k}{\langle \rho_k \rangle^T \lambda_k} < U_b \quad (13)$$

where  $U_b$  is upper bound.

## 3. Simulation example

(See Table 2)

### 3.1 First case : XOR problem

(See Tables 3, 4)

### 3.2 Second case : Evaluation of Rosenbrock function value

(See Tables 3, 5)

[Table 2] Input parameters

Simulation	SDBP				QNBP	
	Step Size	Momentum -Rate	Initial Temperature	Type	$U_b$	$s$
1	1.0	0.5	0	d*	1.0	1.0
2	1.0	0.5	0.001	s**	0.5	1.0
3	0.5	0.25	0	d	0.3	1.5
4	0.5	0.25	0.001	s	0.2	2.0
5	0.3	0.15	0	d	0.15	2.0
6	0.3	0.15	0.001	s	0.1	2.0
7	0.2	0.1	0	d		
8	0.2	0.1	0.001	s		

\* : deterministic \*\* : stochastic

[Table 3] Simulation environment in two cases

	Input Nodes	First Hidden Nodes	Output Nodes	Output Function	Error Level	Max. Iteration	Initial Weights	Training Data
1_st Case	2	2	1	sigmoid	1E-12	3000	n*	4
2_nd Case	14	10	7	sigmoid	1E-12	5000	n	20

\* : normalized

[Table 4] Simulation results of 1<sub>st</sub> case, XOR problem

Simulation	SDBP			QNBP (DFP+SR1 AND BFGS+SR1)		
	Run Time	Iterations	Error	Run Time	Iterations	Error
1	0.77	3000	0.00050596986	0.33	500	0.000000000000
2	0.77	3000	0.00051034391	0.27	400	0.000000000000
3	0.77	3000	0.00224905702	0.33	500	0.000000000000
4	0.77	3000	0.00223123949	0.33	500	0.000000000000
5	0.77	3000	0.00853072653	0.33	500	0.000000000000
6	0.77	3000	0.00825091817	0.49	700	0.000000000000
7	0.77	3000	0.09576089606			
8	0.77	3000	0.08420758138			

[Table 5] Simulation results of 2<sub>nd</sub> case, Rosenbrock function evaluation

Simulation	SDBP			QNBP (DFP+SR1 AND BFGS+SR1)		
	Run Time	Iterations	Error	Run Time	Iterations	Error
1	61.48	5000	0.00485982260	133.0	1000	0.000000000001
2	61.48	5000	0.00144454306	292.5	2200	0.000000000001
3	61.48	5000	0.00450409402	412.0	3100	0.000000000001
4	61.48	5000	0.00462669594	465.4	3500	0.000000000001
5	61.48	5000	0.00764988971	531.21	4000	0.000000000001
6	61.48	5000	0.00752127151	611.76	4600	0.000000000001
7	61.48	5000	0.01256565752			
8	61.48	5000	0.01190128557			

#### 4. Concluding remarks

From simulation results, compared with SDBP, QNBP has three major advantages. These are as follows.

① QNBP is by far general, i.e., SDBP is a special case of QNBP  
 (  $A_k = I \quad \forall k$  )

② QNBP is definitely superior to SDBP in terms of accuracy and rate of convergence.

③ QNBP may surmount the problem of local minima.  
 And

① In large scale problem, QNBP is time-consuming.

② There was no difference between DFP+SR1 and BFGS+SR1 combined algorithms in QNBP.

#### References

- [1] W.C. Davidon, "Variable Metric Method for Minimization", Argonne National Laboratory Report ANL-5990, Argonne, Illinois, 1959.
- [2] R. Fletcher and M.J.D. Powell, "A Rapidly Convergent Descent Method for Minimization", Computer Journal, 6, 163-168, 1963.
- [3] C.G. Broyden, "The Convergence of a Class of Double-Rank Minimization Algorithms", Journal Inst. Maths. Applics., 6, 76-90, 1970.
- [4] R. Fletcher, "A New Approach to Variable Metric Algorithm",

Computer Journal, 13, 392-399, 1970.

- [5] D. Goldfarb, "A Family of Variable Metric Methods Derived by Variational Means", Mathematics of Computations 24, 23-26, 1970.
- [6] D.F. Shanno, "Conditioning of Quasi-Newton Methods for Function Minimization", Mathematics of Computations 24, 647-657, 1970.
- [7] W.C. Davidon, "Variance Algorithm for Minimization", Computer Journal, 10, 406-410, 1968.