

# ADAPTIVE MULTIPLE OPTIMAL LEARNING FACTORS FOR NEURAL NETWORK TRAINING

---

Jeshwanth Challagundla  
Supervising Prof: Dr. Michael T. Manry  
April 10 2015

# Outline

- I. Multilayer perceptron review
- II. Review of neural network training algorithms
- III. Goals
- IV. Proposed algorithm
- V. Simulations
- VI. Conclusions

# I. MLP Review

- ❑ Widely used for function approximation and pattern recognition
  - In function approximation a function that maps between input and output is approximated
  - In pattern recognition a received signal is assigned to one of a prescribed number of classes
- ❑ Known for its special properties of approximating Bayes discriminants and universal approximation
- ❑ Function approximation is the focus of this thesis

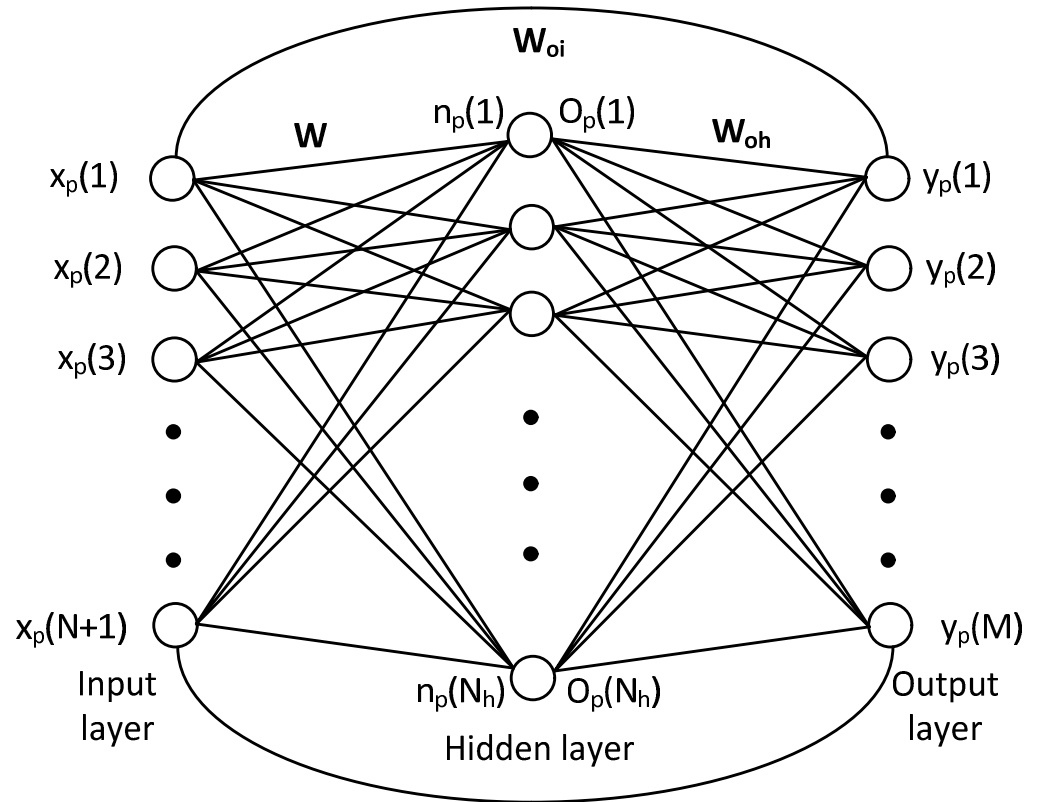
$\mathbf{W}$   $\rightarrow$  Input weights  
 $\mathbf{W}_{oi}$   $\rightarrow$  Bypass weights  
 $\mathbf{W}_{oh}$   $\rightarrow$  Output to hidden weights

$$\mathbf{n}_p = \mathbf{W} \cdot \mathbf{x}_p$$

$$O_p(k) = \frac{1}{1 + e^{-n_p(k)}}$$

$$\mathbf{y}_p = \mathbf{W}_{oi} \cdot \mathbf{x}_p + \mathbf{W}_{oh} \cdot \mathbf{O}_p$$

$$E = \frac{1}{N_v} \sum_{p=1}^{N_v} \sum_{i=1}^M [t_p(i) - y_p(i)]^2$$



## II. Review of NN training algorithms

### □ Conjugate Gradient

Let,  $\mathbf{w} = \text{vec}\{\mathbf{W}, \mathbf{W}_{oi}, \mathbf{W}_{oh}\}$  ,  $\mathbf{g} = \text{vec}\{\mathbf{G}, \mathbf{G}_{oi}, \mathbf{G}_{oh}\}$

$$\mathbf{p} = \text{vec}\{\mathbf{P}, \mathbf{P}_{oi}, \mathbf{P}_{oh}\}$$

$$\mathbf{p}_{k+1} = \mathbf{g}_k + \mathbf{B}_1 \cdot \mathbf{p}_k$$

$$\mathbf{B}_1 = \frac{\mathbf{g}_{k+1}^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k} , \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{z} \cdot \mathbf{p}_k$$

- Minimizes quadratic error functions of  $n$  variables in  $n$  steps
- No Hessian computation is required

## □ Levenberg Marquardt

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{e}$$

$$\mathbf{e} = [\mathbf{H}' + \lambda \mathbf{I}]^{-1} \mathbf{g}'$$

$$h'(m, n) = \frac{\partial^2 E}{\partial w(m) \partial w(n)} \quad g'(m) = \frac{-\partial E}{\partial w(m)}$$

- By changing  $\lambda$  the algorithm interpolates between steepest descent and Gauss-Newton methods
- Computationally very expensive
- Mostly used for small networks

# Newton-related improved training

Newton-related algorithms usually converge faster per iteration

## □ Output Weight Optimization

Augmented input:

$$\mathbf{R} = \frac{1}{N_v} \sum_{p=1}^{N_v} \mathbf{X}_p \mathbf{X}_p^T$$

$$\mathbf{C} = \frac{1}{N_v} \sum_{p=1}^{N_v} \mathbf{X}_p \mathbf{t}_p^T$$

$$\mathbf{X}_p = \begin{bmatrix} \mathbf{x}_p \\ \mathbf{O}_p \end{bmatrix}$$

$$\mathbf{W}_o^T = \mathbf{R}^{-1} \cdot \mathbf{C}$$

- Equivalent to Newton's method for output weights
- Computationally less expensive

## ❑ OWO-MOLF

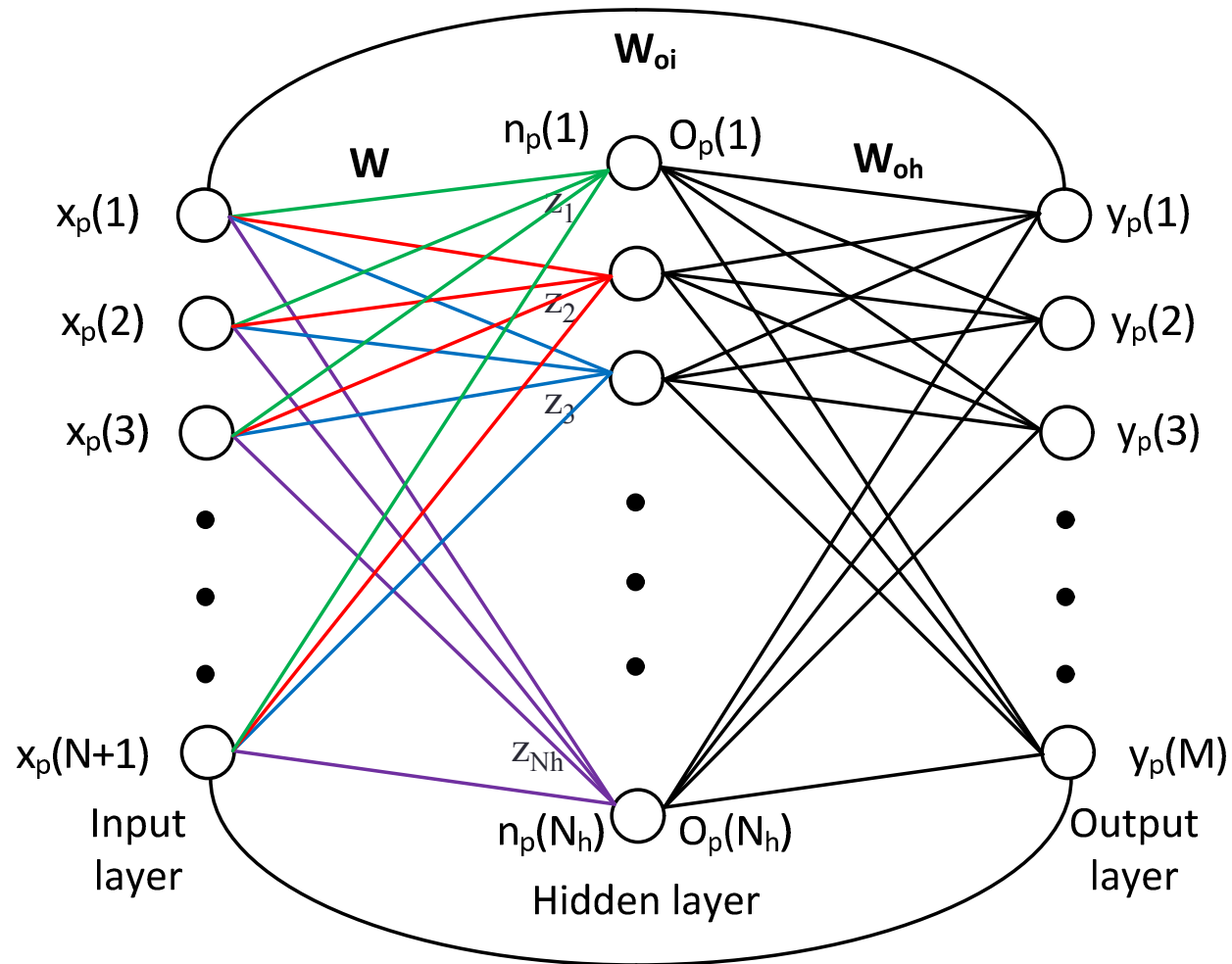
In each iteration:

- Input weights are trained using
  - Learning factor vector  $\mathbf{z}$  of size  $N_h \times 1$
  - Negative gradient vector  $\mathbf{g}$
- Output weights are trained using OWO
- $z_k$ ,  $k^{\text{th}}$  element of vector  $\mathbf{z}$  used to update all the weights connected to hidden unit  $k$

$$y_p(i) = \sum_{n=1}^{N+1} w_{oi}(i, n) x_p(n) + \sum_{k=1}^{N_h} w_{oh}(i, k) f\left(\sum_{n=1}^{N+1} (w(k, n) + z_k g(k, n)) x_p(n)\right)$$



# MOLF demonstration



- Learning factor vector  $\mathbf{z}$  is computed as

$$\mathbf{z} = \mathbf{H}_{\text{molf}}^{-1} \mathbf{g}_{\text{molf}}$$

- Element of matrix  $\mathbf{H}_{\text{molf}}$  ,

$$h_{\text{molf}}(k, j) = \frac{\partial^2 E}{\partial z_k \partial z_j}$$

- $\mathbf{g}_{\text{molf}} = \left[ -\partial E / \partial z_1, -\partial E / \partial z_2, \dots, -\partial E / \partial z_{N_h} \right]^T$
- OWO-MOLF trains much faster than BP

## □ OWO-Newton

In each iteration:

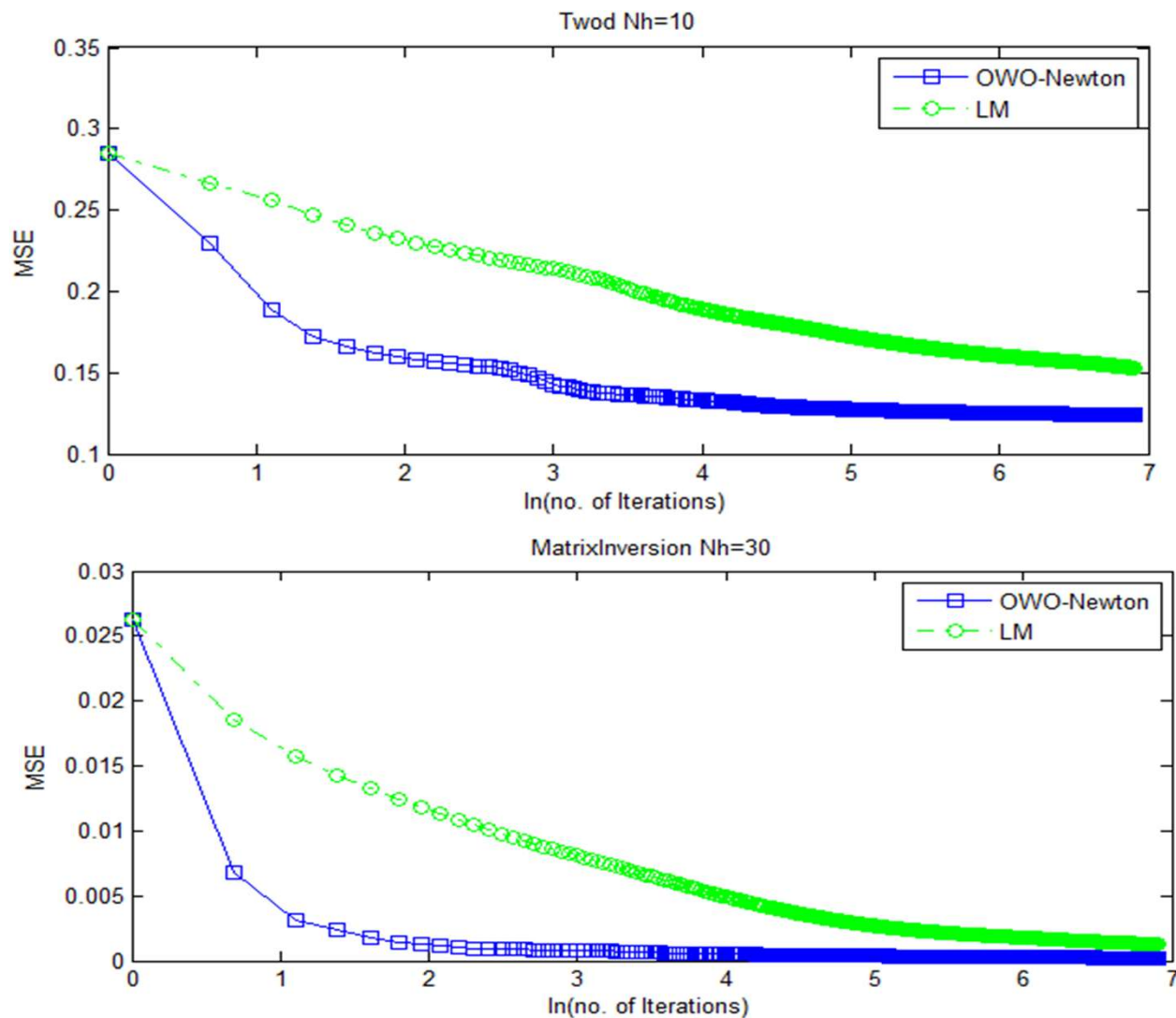
- Newton's method for input weights

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{e} \qquad \mathbf{w} = \text{vec}\{\mathbf{W}\}$$

$$\mathbf{e} = \mathbf{H}^{-1} \cdot \mathbf{g}$$

- $\mathbf{H}$  → Hessian of error with respect to input weights in the network
- $\mathbf{g}$  → negative input weight gradient vector
- Output weights are solved using OWO

- OWO-Newton has simpler and faster code than LM for similar performance



## ❑ Problems

- Little theoretical justification for using small second order modules in first order algorithms
- Second order training is computationally very expensive
- Error decrease per iteration is larger in OWO-Newton, number of multiplies per iteration is less for OWO-MOLF

## □ Lemmas

➤ Lemma 1: Assume for quadratic  $E(\mathbf{w})$ ,  $\mathbf{w}$  is divided into  $k$  partitions  $\mathbf{w}_k$ . If we minimize  $E$  with respect to  $\mathbf{z}$  producing an error  $E_k = E(\mathbf{w}_1 + z_1 \mathbf{g}_1, \mathbf{w}_2 + z_2 \mathbf{g}_2, \dots, \mathbf{w}_k + z_k \mathbf{g}_k)$  and  $k$  increases by splitting an existing partition, then  $E_{k+1} \leq E_k$

## ➤ Implications

- More learning factors means more error decrease
- Not guaranteed since  $E(\mathbf{w})$  is not quadratic

### III. Goals

- ❑ Develop a new algorithm that adapts between OWO-MOLF and OWO-Newton
- ❑ Maximize error decrease per multiply
- ❑ Simulations

## IV. Adaptive Multiple Optimal Learning Factors

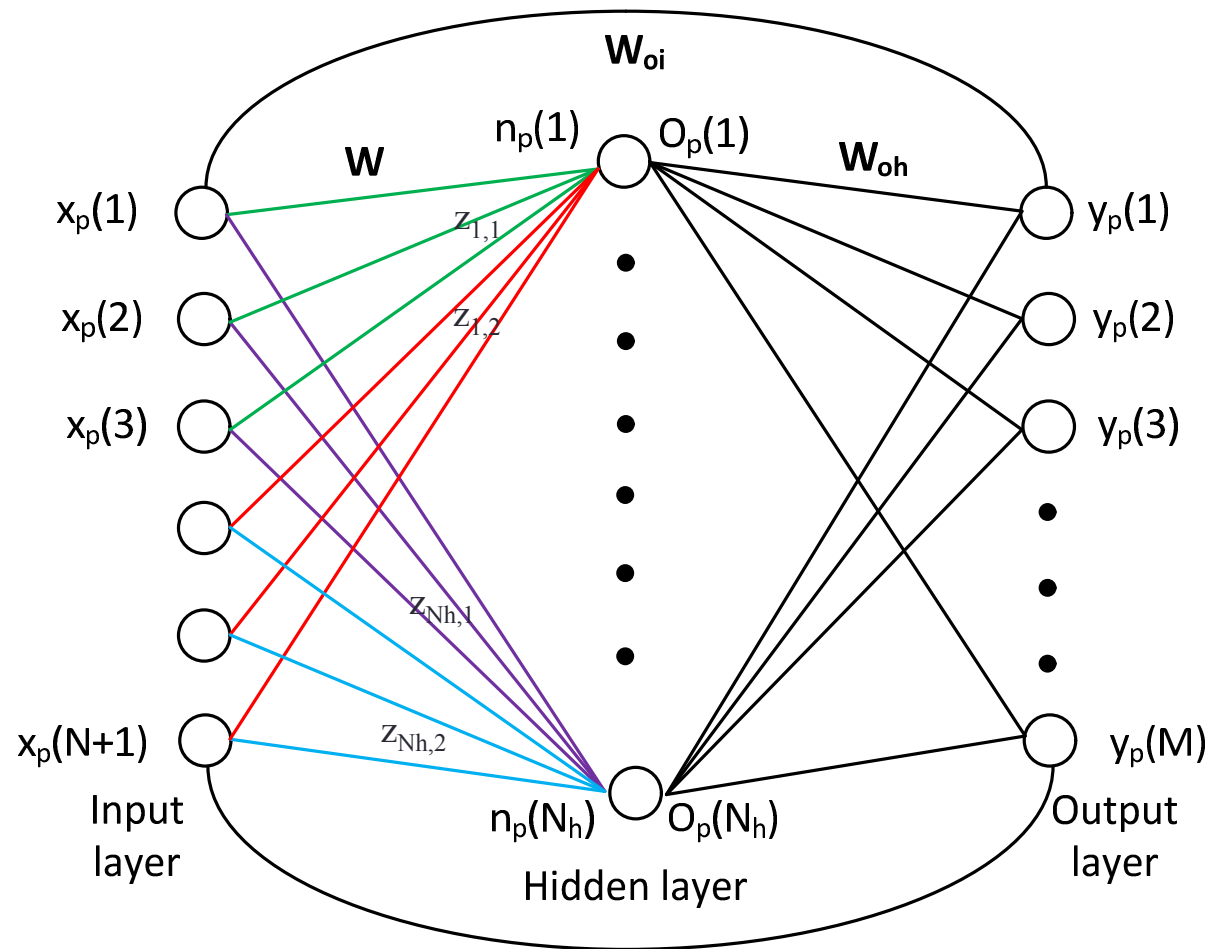
In each iteration:

- Input weights are trained using
  - Learning factor vector  $\mathbf{z}$
  - Negative gradient vector  $\mathbf{g}$
- Output weights are trained using OWO
- Size of  $\mathbf{z}$  can vary between  $N_h$  and  $N_h \times (N+1)$  in each iteration

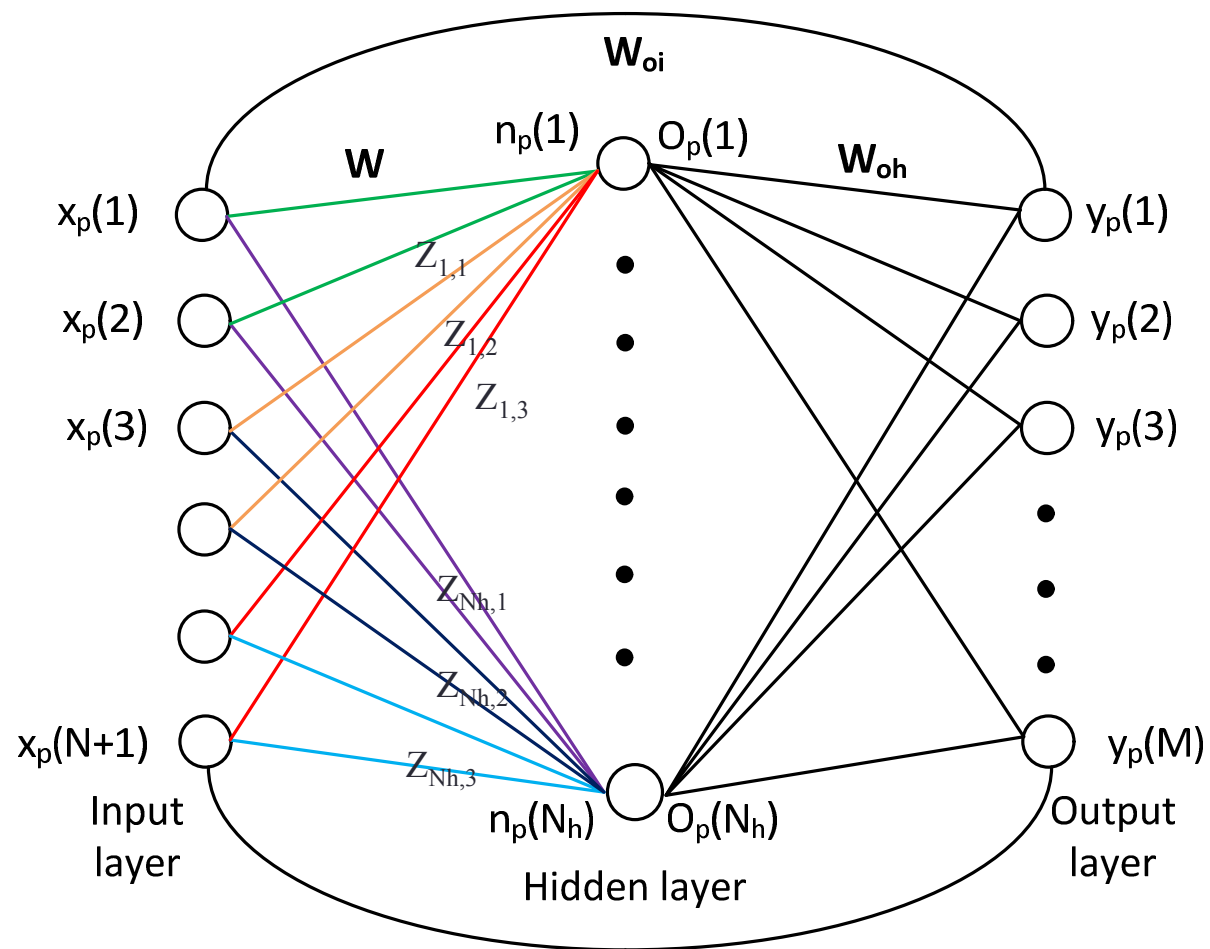


- This algorithm adapts between OWO-MOLF and OWO-Newton
- The input weights connected to each hidden unit are grouped into  $N_g$  groups
- An unique learning factor is computed for each of this groups
- $z_{k,C}$  is an element of  $\mathbf{z}$  used to update all the input weights that belong to the group  $C$  of hidden unit  $k$ .

# Adaptive MOLF with $N_g=2$



# Adaptive MOLF with $N_g=3$



## □ Grouping of input weights

- Grouping is done based on the curvature  $\mathbf{H}_w$  of the error function calculated with respect to input weights
- The elements of  $\mathbf{H}_w$  can be calculated as,

$$h_w(k, n) = \frac{\partial^2 E}{\partial w(k, n)^2} = \frac{2}{N_v} \sum_{i=1}^M w_{oh}(i, k)^2 \sum_{p=1}^{N_v} f'(n_p(k))^2 x_p(n)^2$$

□ How is the number of groups changed

- Goal is to vary number of groups so that the error change per multiply is maximized
- Number of groups per hidden unit increases as the error change per multiply (EPM) increases and vice versa.

$$\text{EPM}(i_t) = \frac{E(i_t - 1) - E(i_t)}{M(i_t)}$$

- $M(i_t)$  stands for number of multiplies in iteration  $i_t$
- $\text{EPM}(i_t)$  is the error per multiply in iteration  $i_t$
- $E(i_t)$  is the error computer in iteration  $i_t$
- $i_t$  is the current iteration number.

## □ Derivation of adaptive MOLF algorithm

$$y_p(i) = \sum_{n=1}^{N+1} x_p(n) w_{oi}(k, n) + \sum_{k=1}^{N_h} w_{oh}(i, k) f \left( \sum_{C=1}^{N_g} \sum_{a=R(C-1)+1}^{R(C)} x_p(i_k(a)) [w(k, i_k(a)) + z_{k,C} g(k, i_k(a))] \right)$$

$$R(C) = C \times Gs(C)$$

$$Gs(0) = 0$$

- C is the group index
- **Gs** is an array of  $N_g \times 1$  elements, these elements contains the sizes of all group the belong to a hidden unit
- **I<sub>k</sub>** contains input indices ordered in such a way that the index n of input to which weight  $w(k, n)$  with higher curvatures are connected will come first.
- **I<sub>k</sub>** =  $[n_1, n_2, n_3, \dots, n_{N+1}]$  where  $n_1, n_2, n_3, \dots, n_{N+1}$  are input indices such that  $h_w(k, n_1) \geq h_w(k, n_2) \geq h_w(k, n_3) \dots \geq h_w(k, n_{N+1})$

## □ Hessian and negative gradient elements

$$\mathbf{g}_{\text{Amolf}}(\mathbf{k}, \mathbf{C}) = \frac{-\partial E}{\partial \mathbf{z}_{\mathbf{k}, \mathbf{C}}} = \frac{2}{N_v} \sum_{p=1}^{N_v} \sum_{i=1}^M (t_p(i) - y_p(i)) \frac{\partial y_p(i)}{\partial \mathbf{z}_{\mathbf{k}, \mathbf{C}}}$$

$$\frac{\partial y_p(i)}{\partial \mathbf{z}_{\mathbf{k}, \mathbf{C}}} = w_{\text{oh}}(i, k) f(n_p(k)) \Delta n_p(k, \mathbf{C})$$

$$\Delta n_p(k, \mathbf{C}) = \sum_{a=R(\mathbf{C}-1)+1}^{R(\mathbf{C})} \mathbf{x}_p(i_k(a)) g(k, i_k(a))$$

$$h_{\text{Amolf}}(\mathbf{k}, \mathbf{C}_1, \mathbf{j}, \mathbf{C}_2) = \frac{2}{N_v} \sum_{i=1}^M w_{\text{oh}}(i, k) w_{\text{oh}}(i, j) \sum_{p=1}^{N_v} f(n_p(k)) f(n_p(j)) \Delta n_p(k, \mathbf{C}_1) \Delta n_p(j, \mathbf{C}_2)$$

- The elements of the 2 dimensional Hessian  $\mathbf{H}_{\text{Amolf}}$  are found from  $h_{\text{Amolf}}(k, C_1, j, C_2)$  as

$$h_{\text{Amolf}}((k-1)N_g + C_1, (j-1)N_g + C_2) = h_{\text{Amolf}}(k, C_1, j, C_2)$$

- Elements of negative gradient column vector  $\mathbf{g}_{\text{Amolf}}$  are calculated from  $g_{\text{Amolf}}(k, C)$  as,

$$g_{\text{Amolf}}((k-1)N_g + C) = g_{\text{Amolf}}(k, C)$$

$$\mathbf{z} = \mathbf{H}_{\text{Amolf}}^{-1} \cdot \mathbf{g}_{\text{Amolf}}$$



□ The input weights  $\mathbf{W}$  are updated as follows,

$$w(k, i_k(a)) = w(k, i_k(a)) + z_{k,C} g(k, i_k(a))$$

where,

$$[(C-1) \times Gs(C-1)] + 1 \leq a \leq [C \times Gs(C)]$$

- Lemma 2:  $E - E_{\text{MOLF}}$  and  $E - E_{\text{aMOLF}}$  denote the error decrease due to the Newton steps of OWO-MOLF and adaptive MOLF respectively, then  $E - E_{\text{MOLF}} \leq E - E_{\text{aMOLF}}$
- Lemma 3: OWO-Newton is a limiting case of the adaptive MOLF algorithm as the  $k$  groups of adaptive MOLF are split until  $k = N_h \cdot (N + 1)$

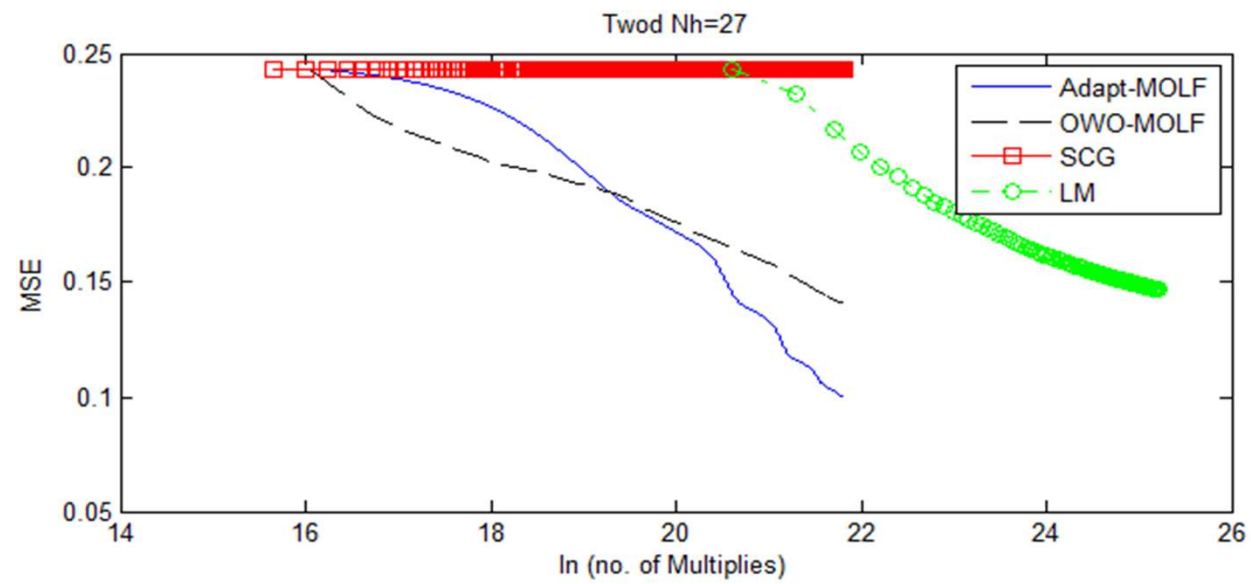
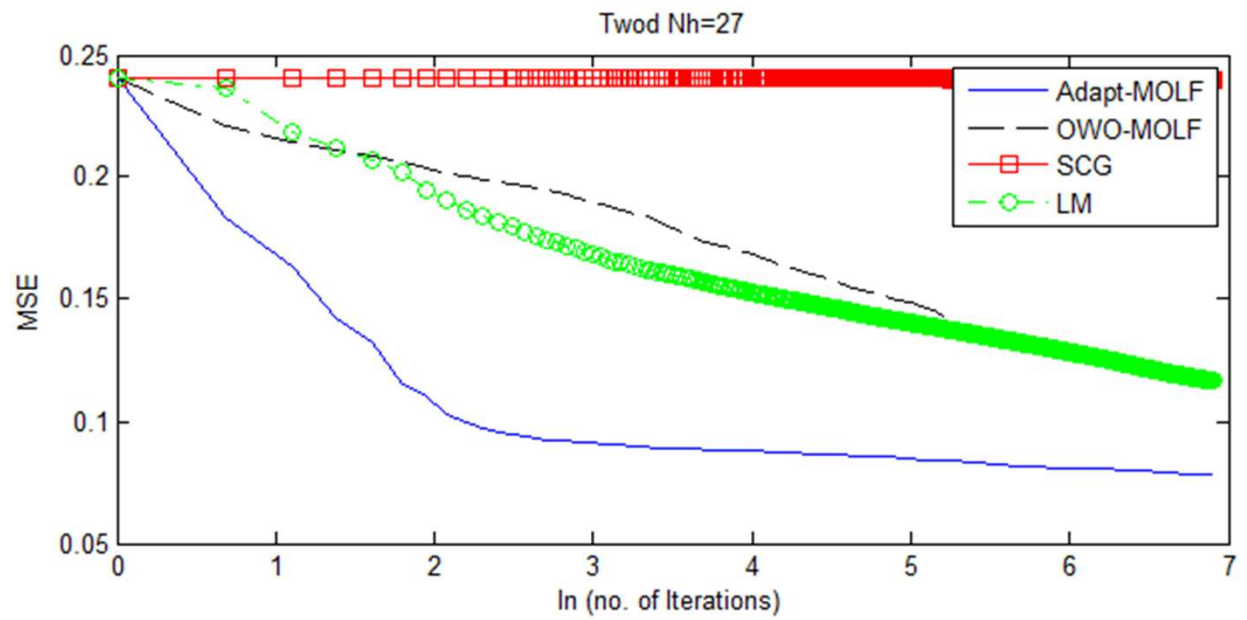
- Adaptive MOLF Hessian and gradients can be generated from Gauss-Newton Hessian and gradients as follows

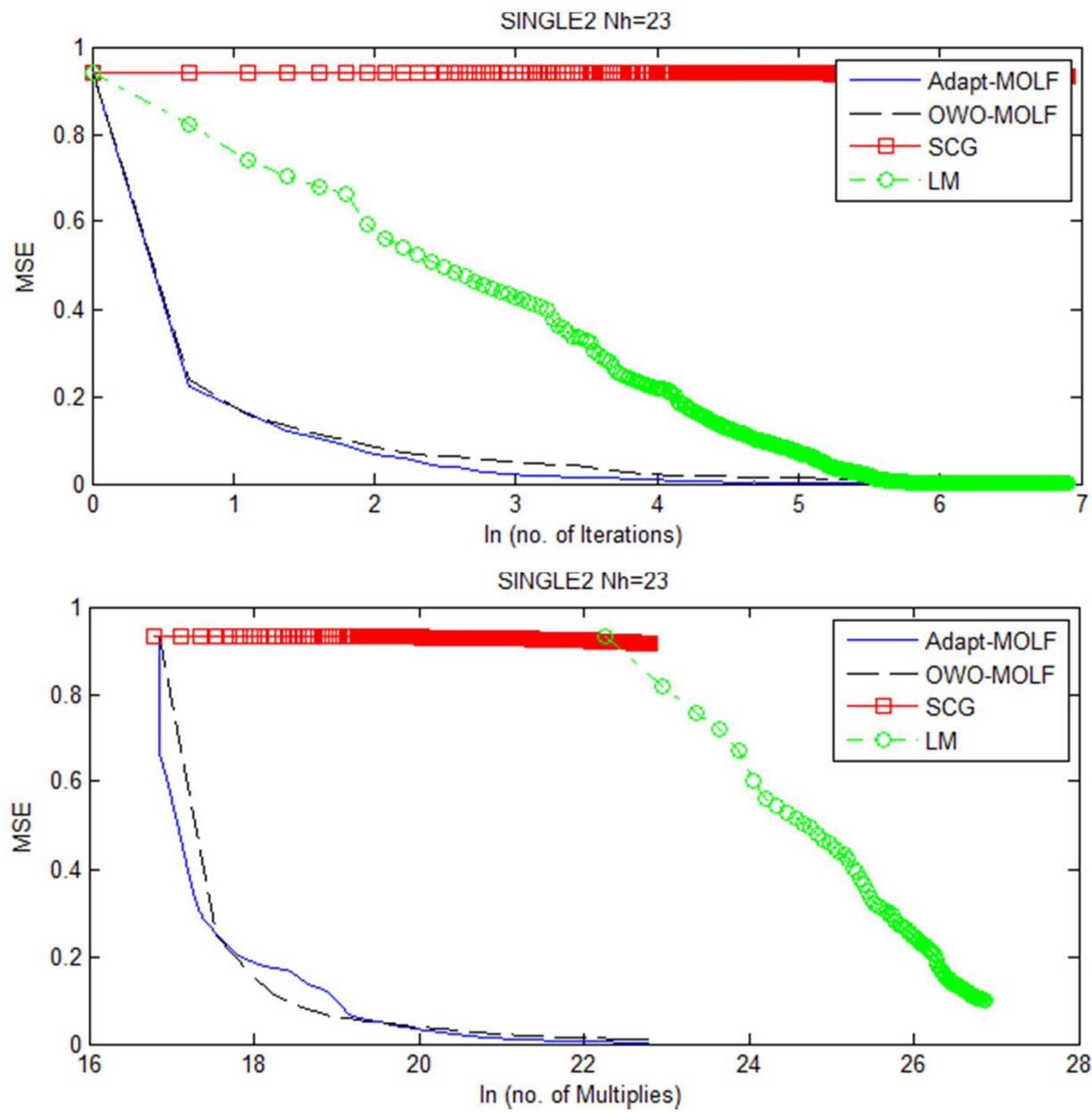
$$g_{\text{Amolf}}(k, C) = \sum_{a=R(C-1)+1}^{R(C)} g(k, i_k(a))^2$$

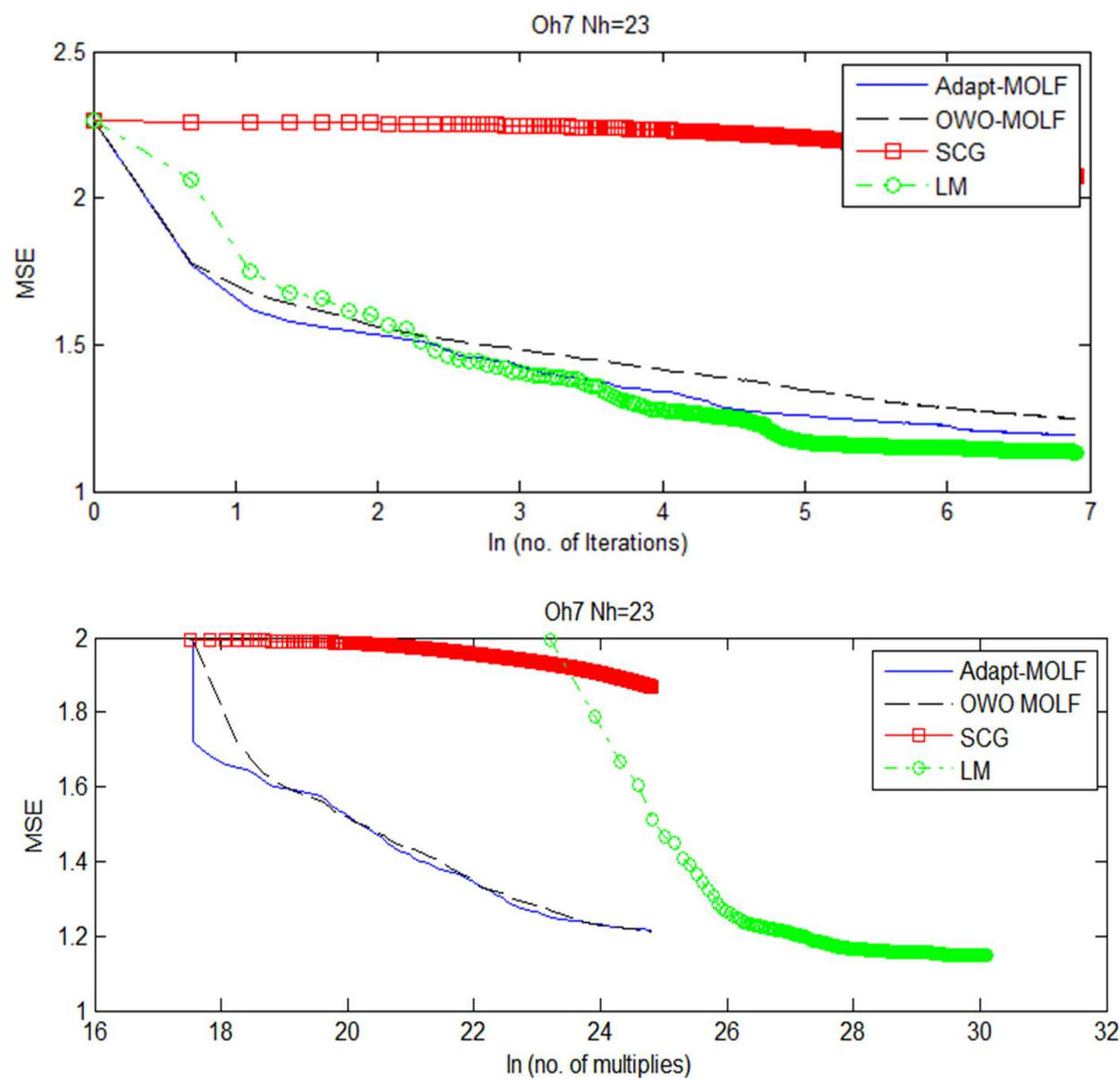
$$h_{\text{Amolf}}(k, C_1, j, C_2) = \sum_{a=R(C_1-1)+1}^{R(C_1)} \sum_{b=R(C_2-1)+1}^{R(C_2)} h(k, i_k(a), j, i_j(b)) g(k, i_k(a)) g(j, i_j(b))$$

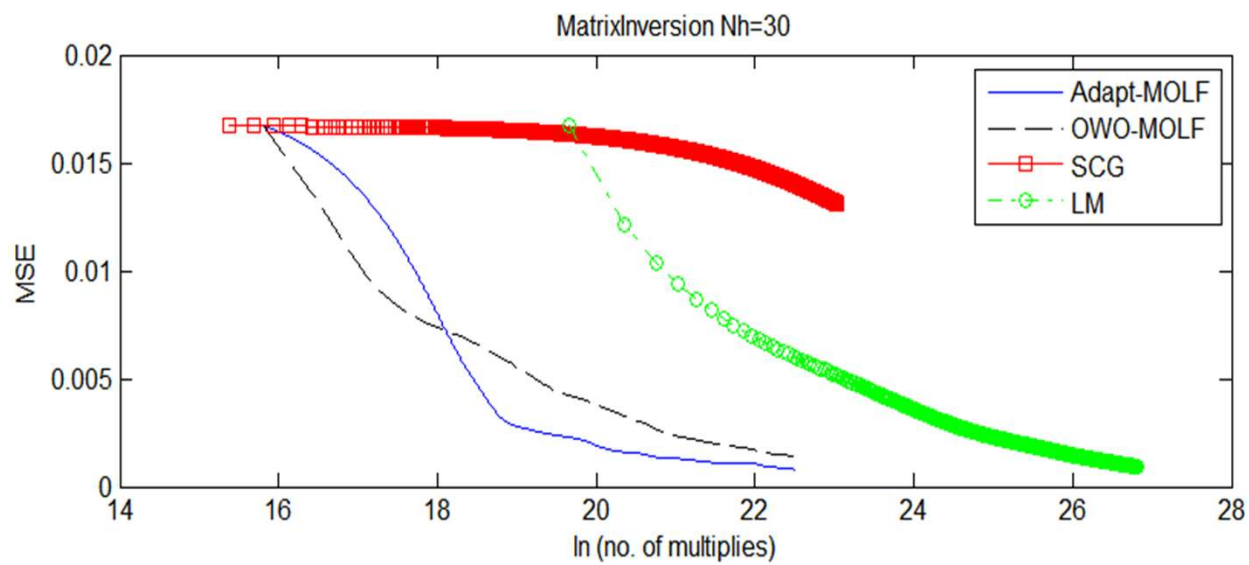
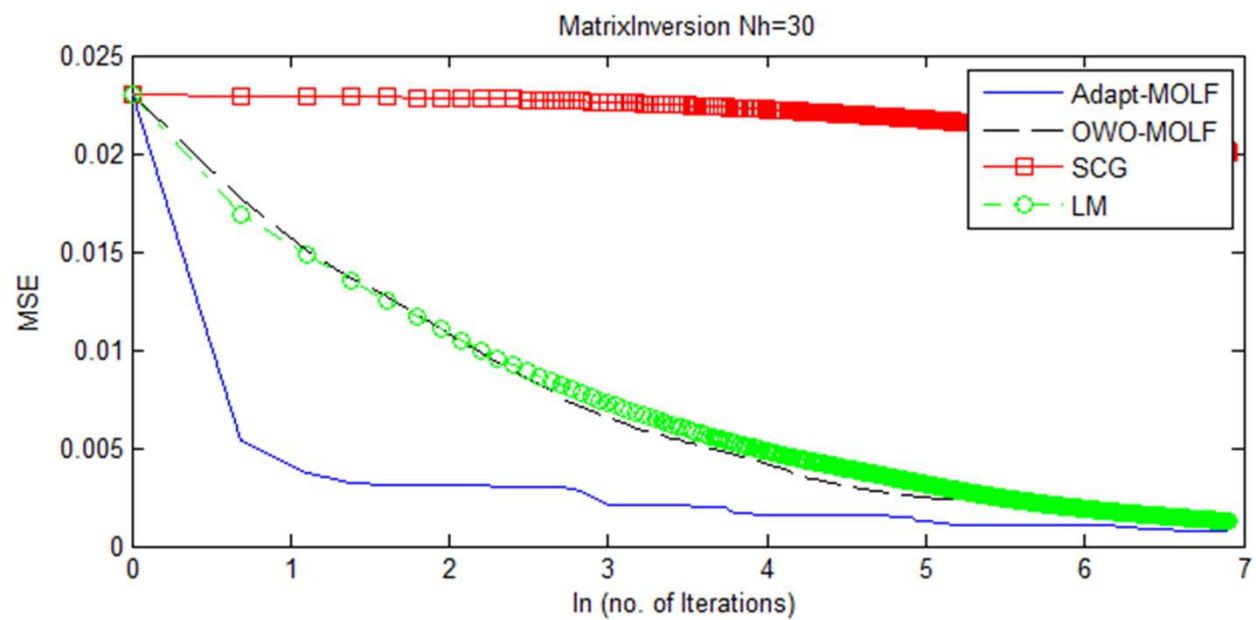
## V. Simulations

- Training is done on the entire data set 10 times with 10 different initial networks. The average MSE from this 10-fold training is shown in the plots below.











## □ K-fold validation and testing

Data Set		Adaptive MOLF	OWO-MOLF	SCG	LM
Twod.tra	$E_{\text{trn}}$	0.0888	0.1554	1.0985	0.2038
	$E_{\text{tst}}$	0.1172	0.1731	1.0945	0.2205
Single2.tra	$E_{\text{trn}}$	0.0042	0.0151	3.5719	0.0083
	$E_{\text{tst}}$	0.2319	0.1689	3.6418	0.0178
Mattrn.tra	$E_{\text{trn}}$	0.0011	0.0027	4.2400	0.0022
	$E_{\text{tst}}$	0.0013	0.0032	4.3359	0.0027
Oh7.tra	$E_{\text{trn}}$	1.2507	1.3205	4.1500	1.1602
	$E_{\text{tst}}$	1.4738	1.4875	4.1991	1.4373

## VI. Conclusions

- ❑ The number of learning factors needed by a training algorithm is addressed and adaptively optimized
- ❑ The adaptive MOLF algorithm is superior to the OWO-MOLF algorithm in terms of error decrease per iteration and often in terms of error decrease per multiply
- ❑ The proposed algorithm interpolates between OWO-MOLF and OWO-Newton



Thank you



Questions