

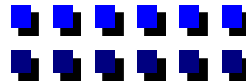


# Lecture 6: Variable Metric (Secant, Quasi-Newton) Methods, Quadratic Termination, Practicalities, Incremental Methods for NN

**Prof. Krishna R. Pattipati**  
**Dept. of Electrical and Computer Engineering**  
**University of Connecticut**  
**Contact: [krishna@engr.uconn.edu](mailto:krishna@engr.uconn.edu) (860) 486-2890**

***ECE 6437***  
***Computational Methods for Optimization***

***Fall 2009***  
***October 6, 2009***





## Outline of Lecture 6

- Quasi-Newton (QN) Methods
- Motivation from Quadratic Functions
- Square-root Implementation
  - Davidon-Fletcher-Powell (DFP) & Broyden-Fletcher-Goldfarb-Shanno (BFGS) versions
  - Relation to Kalman Filtering with Perfect Observations
- Properties of QN Methods
  - Quadratic Termination
  - Convergence Properties
  - Scaling
- Incremental Methods for Neural Networks



# Quasi-Newton Methods - 1

- We want to look at a class of algorithms of the form:

$$\underline{x}_{k+1} = \underline{x}_k - \alpha_k H_k \underline{g}_k, \quad \underline{g}_k = \nabla f(\underline{x}_k)$$

$$H_k \sim [\nabla^2 f(\underline{x}_k)]^{-1} \Rightarrow \text{Newton is fast, but requires Hessian}$$

$$H_k \sim I \Rightarrow \text{SD}$$

$$\underline{d}_k = -\underline{g}_k + \frac{(\underline{g}_k - \underline{g}_{k-1})^T \underline{g}_k}{(\underline{g}_k - \underline{g}_{k-1})^T \underline{d}_{k-1}} \underline{d}_{k-1} = - \left[ I - \frac{\underline{d}_{k-1} \underline{q}_{k-1}^T}{\underline{q}_{k-1}^T \underline{d}_{k-1}} \right] \underline{g}_k = -H_k \underline{g}_k$$

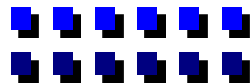
$$\underline{q}_{k-1} = \underline{g}_k - \underline{g}_{k-1}$$

$H_k$  not symmetric, but can be made symmetric. See Shanno in *M. of OR* (1978)

- **Key Question:** Can we build up  $H_k \ni H_0 \rightarrow H_1 \rightarrow \dots \rightarrow H_k \rightarrow [\nabla^2 f(\underline{x}_k)]^{-1}$  without explicitly computing the Hessian?

Essentially an extension of **secant** approximation to second derivatives:

$$\frac{d^2 f}{dx^2} \cong \frac{(df/dx)|_{x_1} - (df/dx)|_{x_2}}{x_1 - x_2} \text{ to higher dimensions.}$$





# Quasi-Newton Methods - 2

- Sometimes we may want an approximation to Hessian

$$B_0 \rightarrow B_1 \rightarrow \dots \rightarrow B_k \rightarrow \nabla^2 f(\underline{x}_k)$$

- Let us recall the requirements on  $H_k$  and  $B_k$ :

1.  $H_k$  is PD (or)  $H_k > 0$

$B_k > 0$

2.  $H_k$  symmetric  $\Rightarrow H_k = H_k^T$   $\longleftrightarrow$

$B_k = B_k^T$

3. Finite termination on quadratics

quadratic termination

4.  $H_k \sim [\nabla^2 f(\underline{x}_k)]^{-1}$

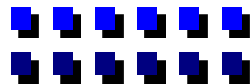
$B_k \sim \nabla^2 f(\underline{x}_k)$

- Process:  $\underline{x}_k \rightarrow \underline{g}_k \xrightarrow{\Delta} H_k \rightarrow \underline{d}_k \rightarrow \alpha_k \rightarrow \underline{x}_{k+1} \rightarrow \underline{g}_{k+1} \xrightarrow{\Delta} H_{k+1}$

- Motivation via quadratic functions:  $f(\underline{x}) = \frac{1}{2}(\underline{x} - \underline{x}^*)^T Q(\underline{x} - \underline{x}^*)$

$$\underline{g}_k = Q(\underline{x} - \underline{x}^*), \quad \underline{g}_{k+1} = Q(\underline{x}_k - \underline{x}^* + \underline{x}_{k+1} - \underline{x}_k) = \underline{g}_k + Q\underline{p}_k$$

where  $\underline{p}_k = \underline{x}_{k+1} - \underline{x}_k = \alpha_k \underline{d}_k$





# Motivation via Quadratic Functions

## Key Secant Relations:

Look for relations of the form:

$$Q \underline{p}_k = \underline{q}_k = \underline{g}_{k+1} - \underline{g}_k \quad \Rightarrow \quad B_{k+1} \underline{p}_k = \underline{q}_k$$

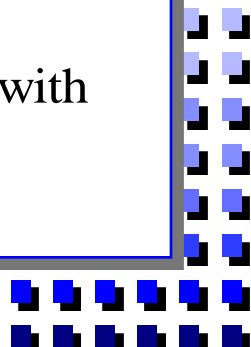
$$\text{(or)} \quad Q^{-1} \underline{q}_k = \underline{p}_k \quad \Rightarrow \quad H_{k+1} \underline{q}_k = \underline{p}_k$$

$$Q \begin{bmatrix} \underline{p}_0 & \underline{p}_1 & \cdots & \underline{p}_{n-1} \end{bmatrix} = \begin{bmatrix} \underline{q}_0 & \underline{q}_1 & \cdots & \underline{q}_{n-1} \end{bmatrix}$$

$$\Rightarrow Q \sim B = \begin{bmatrix} \underline{q}_0 & \underline{q}_1 & \cdots & \underline{q}_{n-1} \end{bmatrix} \begin{bmatrix} \underline{p}_0 & \underline{p}_1 & \cdots & \underline{p}_{n-1} \end{bmatrix}^{-1}$$

$$\Rightarrow Q^{-1} \sim H = \begin{bmatrix} \underline{p}_0 & \underline{p}_1 & \cdots & \underline{p}_{n-1} \end{bmatrix} \begin{bmatrix} \underline{q}_0 & \underline{q}_1 & \cdots & \underline{q}_{n-1} \end{bmatrix}^{-1}$$

- We can approximate  $Q$  or  $Q^{-1}$  via  $\{\underline{p}_i\}$  and  $\{\underline{q}_i\}$ , but it may not be symmetric.
- Can we do this recursively?  $H_k \rightarrow H_{k+1}$  (or)  $B_k \rightarrow B_{k+1}$  starting with  $H_0$  or  $B_0$ . Typically  $H_0 = B_0 = I$ .





# Key Ideas of Quasi-Newton Methods - 1

## □ Problem:

### Inverse Hessian Approximation

Given:  $\underline{p}_k = \underline{x}_{k+1} - \underline{x}_k$

$$\underline{q}_k = \underline{g}_{k+1} - \underline{g}_k$$

$$H_k > 0 \ \& \ H_k = H_k^T$$

Find:  $H_{k+1} = H_{k+1}^T > 0$

$$\ni H_{k+1} \underline{q}_k = \underline{p}_k$$

### Hessian Approximation

Given:  $\underline{p}_k = \underline{x}_{k+1} - \underline{x}_k$

$$\underline{q}_k = \underline{g}_{k+1} - \underline{g}_k$$

$$B_k > 0 \ \& \ B_k = B_k^T$$

Find:  $B_{k+1} = B_{k+1}^T > 0$

$$\ni B_{k+1} \underline{p}_k = \underline{q}_k$$

(DUALS)

Secant (or) quasi-Newton conditions

## □ Alternate Problem:

Given:  $\underline{p}_k, \underline{q}_k, L_k$

Find:  $L_{k+1} \ni H_{k+1} = L_{k+1} L_{k+1}^T$

$$\ \& \ L_{k+1} L_{k+1}^T \underline{q}_k = \underline{p}_k$$

Given:  $\underline{p}_k, \underline{q}_k, \hat{L}_k$

Find:  $\hat{L}_{k+1} \ni B_{k+1} = \hat{L}_{k+1} \hat{L}_{k+1}^T$

$$\ \& \ \hat{L}_{k+1} \hat{L}_{k+1}^T \underline{p}_k = \underline{q}_k$$



# Key Ideas of Quasi-Newton Methods - 2

- What are the conditions that guarantee positive definiteness of  $H_{k+1}$ ?

Need  $\underline{p}_k^T \underline{q}_k > 0$

- Scalar Case:  $h_{k+1} = L_{k+1}^2 = \frac{p_k}{q_k}$  need  $p_k, q_k \neq 0$  &  $p_k, q_k$  same sign.

$$H_{k+1} = H_{k+1}^T > 0 \iff \underline{p}_k^T \underline{q}_k > 0$$

- Suppose we have  $L_{k+1} \ni L_{k+1} L_{k+1}^T \underline{q}_k = \underline{p}_k$

$$\text{Let } \underline{v}_k = L_{k+1}^T \underline{q}_k \implies 0 < \underline{v}_k^T \underline{v}_k = \underline{q}_k^T L_{k+1} L_{k+1}^T \underline{q}_k = \underline{q}_k^T \underline{p}_k$$

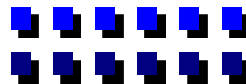
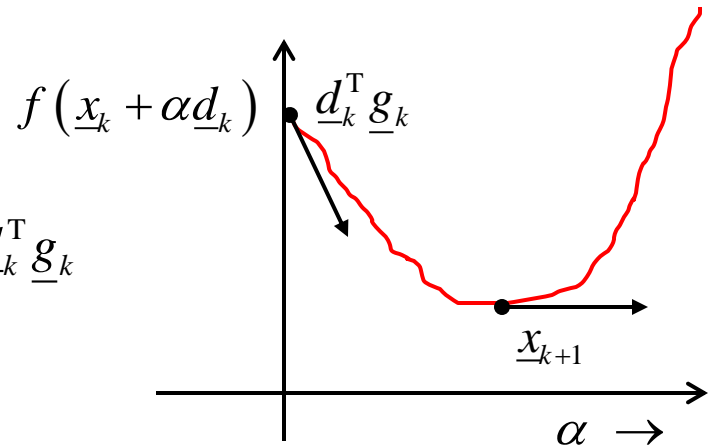
- What does it mean?

$$\underline{p}_k^T \underline{q}_k > 0 \implies \underline{p}_k^T \underline{g}_{k+1} > \underline{p}_k^T \underline{g}_k$$

$$\text{since } \underline{p}_k = \alpha_k \underline{d}_k \text{ \& } \alpha_k > 0 \implies \underline{d}_k^T \underline{g}_{k+1} > \underline{d}_k^T \underline{g}_k$$

$$\underline{d}_k^T \underline{g}_{k+1} = 0 \text{ for exact line searches}$$

$$\underline{d}_k^T \underline{g}_k < 0 \text{ descent direction}$$





# Recursive Update Problem - 1

- Suppose  $\underline{p}_k^T \underline{q}_k > 0$ . How do we construct  $H_{k+1} (L_{k+1})$  from  $H_k (L_k)$ ,  $\underline{p}_k$ , and  $\underline{q}_k$  ?

We will construct  $L_{k+1}$  in two steps:

1) Construct  $S_{k+1} \ni S_{k+1} S_{k+1}^T = H_{k+1}$ .  $S_{k+1}$  not lower triangular.

2) Use QR factorization to construct  $L_{k+1} \ni S_{k+1} = L_{k+1} T_{k+1}^T$

$L_{k+1} \sim$  Lower triangular

$T_{k+1} \sim$  orthogonal matrix  $\Rightarrow T_{k+1}^{-1} = T_{k+1}^T$

**Note:** We do not need to store  $T_{k+1}$  since  $S_{k+1} S_{k+1}^T = L_{k+1} L_{k+1}^T = H_{k+1}$ .

- Quasi-Newton (or secant) condition  $H_{k+1} \underline{q}_k = \underline{p}_k$  (or)  $S_{k+1} S_{k+1}^T \underline{q}_k = \underline{p}_k$  implies:

$$S_{k+1}^T \underline{q}_k = \underline{u}_k, \quad S_{k+1} \underline{u}_k = \underline{p}_k$$

- So, the problem of finding  $H_{k+1} \ni H_{k+1} \underline{q}_k = \underline{p}_k$  &  $\|H_{k+1} - H_k\|_F$  small is equivalent to:

1a) Find  $S_{k+1} \ni S_{k+1} \underline{u}_k = \underline{p}_k$  &  $\|S_{k+1} - L_k\|_F$  is small.

“Least change characterization”

1b) Once  $S_{k+1}$  is known, find  $\underline{u}_k$  from  $S_{k+1}^T \underline{q}_k = \underline{u}_k$ .

$$\|A\|_F = \sqrt{\text{tr}(A^T A)} = \sqrt{\sum_{i=1}^n \sigma_i^2}$$

$\{\sigma_i\} \sim$  singular values





# Recursive Update Problem - 2

## □ Solution of Problem 1a:

$$\min_{S_{k+1}} \frac{1}{2} \|S_{k+1} - L_k\|_F^2 = \min_{S_{k+1}} \frac{1}{2} \text{tr} \left[ (S_{k+1} - L_k)^T (S_{k+1} - L_k) \right] \text{ subject to } S_{k+1} \underline{u}_k = \underline{p}_k$$

$$L(S_{k+1}, \underline{\lambda}) = \frac{1}{2} \text{tr} \left[ (S_{k+1} - L_k)^T (S_{k+1} - L_k) \right] + \underline{\lambda}^T [S_{k+1} \underline{u}_k - \underline{p}_k]$$

$$\frac{\partial L}{\partial S_{k+1}} = S_{k+1} - L_k + \underline{\lambda} \underline{u}_k^T = 0 \Rightarrow S_{k+1} = L_k - \underline{\lambda} \underline{u}_k^T$$

$$\frac{\partial L}{\partial \underline{\lambda}} = S_{k+1} \underline{u}_k - \underline{p}_k = 0 \Rightarrow L_k \underline{u}_k - \underline{p}_k - \underline{\lambda} (\underline{u}_k^T \underline{u}_k) = 0 \text{ or } \underline{\lambda} = \frac{(L_k \underline{u}_k - \underline{p}_k)}{\underline{u}_k^T \underline{u}_k}$$

$$S_{k+1} = L_k + \frac{(\underline{p}_k - L_k \underline{u}_k) \underline{u}_k^T}{\underline{u}_k^T \underline{u}_k}$$

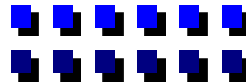
- $S_{k+1}$  is the nearest matrix to  $L_k$  while satisfying  $S_{k+1} \underline{u}_k = \underline{p}_k$ .

## □ Solution of Problem 1b:

$$\text{Need } \underline{u}_k = S_{k+1}^T \underline{q}_k = L_k^T \underline{q}_k + \underline{u}_k \frac{(\underline{p}_k - L_k \underline{u}_k)^T \underline{q}_k}{\underline{u}_k^T \underline{u}_k}$$

$$\left[ 1 - \frac{(\underline{p}_k - L_k \underline{u}_k)^T \underline{q}_k}{\underline{u}_k^T \underline{u}_k} \right] \underline{u}_k = L_k^T \underline{q}_k \Rightarrow \underline{u}_k \parallel L_k^T \underline{q}_k \quad \therefore \underline{u}_k = \Theta_k L_k^T \underline{q}_k$$

$$\begin{aligned} \underline{\lambda}^T S_{k+1} \underline{u}_k &= \text{Tr}(S_{k+1} \underline{u}_k \underline{\lambda}^T) \\ \nabla_A \text{Tr}(AB) &= B^T \\ \Rightarrow \nabla_{S_{k+1}} \text{Tr}(S_{k+1} \underline{u}_k \underline{\lambda}^T) &= \underline{\lambda} \underline{u}_k^T \\ \nabla_{\underline{\lambda}} (\underline{\lambda}^T S_{k+1} \underline{u}_k) &= S_{k+1} \underline{u}_k \end{aligned}$$





# Recursive Update Problem - 3

$$\left[ 1 - \frac{\underline{p}_k^T \underline{q}_k}{\Theta_k^2 \underline{q}_k^T L_k L_k^T \underline{q}_k} + \frac{\underline{q}_k^T L_k^T \underline{u}_k}{\underline{u}_k^T \underline{u}_k} \right] \underline{u}_k = \left[ 1 - \frac{\underline{p}_k^T \underline{q}_k}{\Theta_k^2 \underline{q}_k^T H_k \underline{q}_k} + \frac{1}{\Theta_k} \right] \underline{u}_k = L_k^T \underline{q}_k$$

$$\Rightarrow \Theta_k - \frac{\underline{p}_k^T \underline{q}_k}{\Theta_k \underline{q}_k^T H_k \underline{q}_k} + 1 = 1$$

$\Rightarrow$

$$\Theta_k^2 = \frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}$$

$\Rightarrow$

$$\underline{u}_k = \pm \sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}} \cdot L_k^T \underline{q}_k$$

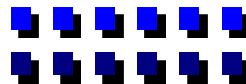
since  $\underline{p}_k^T \underline{q}_k > 0$  sqrt always exists.

- Take the positive sign since this makes  $S_{k+1}$  closer to  $L_k$  when

$$\underline{p}_k \cong H_k \underline{q}_k \Rightarrow \Theta_k = 1 \text{ \& } S_{k+1} \approx L_k \text{ as it should.}$$

- $\underline{u}_k^T \underline{u}_k = \underline{p}_k^T \underline{q}_k > 0$

□ Before we consider the problem of finding  $L_{k+1}$  from  $S_{k+1}$ , let us take a step back and review the results obtained so far.





# Recursive Update Problem - 4

- Given  $L_k$ ,  $\underline{p}_k$ ,  $\underline{q}_k \ni L_k^{-1}$  exists, we found an  $S_{k+1} \ni S_{k+1} S_{k+1}^T \underline{q}_k = \underline{p}_k$ .

The result is:

$$S_{k+1} = L_k + \frac{(\underline{p}_k - L_k \underline{u}_k) \underline{u}_k^T}{\underline{u}_k^T \underline{u}_k}, \quad \underline{u}_k = \Theta_k L_k^T \underline{q}_k$$

$$\Theta_k = \sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}} > 0; \quad H_k = L_k L_k^T$$

- Q:** Is  $S_{k+1}$  invertible?

Suppose it is:

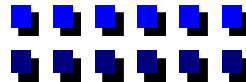
$$S_{k+1}^{-1} = L_k^{-1} - \frac{L_k^{-1} \underline{u}_k (\underline{p}_k - L_k \underline{u}_k)^T L_k^{-1}}{\underline{u}_k^T \underline{u}_k + \underline{u}_k^T L_k^{-1} (\underline{p}_k - L_k \underline{u}_k)}$$

The inverse exists if  $\underline{u}_k^T \underline{u}_k + \underline{u}_k^T L_k^{-1} \underline{p}_k - \underline{u}_k^T \underline{u}_k = \Theta_k \underline{q}_k^T \underline{p}_k \neq 0$ .

Since  $\underline{q}_k^T \underline{p}_k > 0$  &  $\Theta_k > 0$ ,  $S_{k+1}^{-1}$  exists.

- So,  $H_{k+1} = S_{k+1} S_{k+1}^T$  is PD. We can get a recursive expression for  $H_k$  as follows:

$$H_{k+1} = S_{k+1} S_{k+1}^T = \left[ L_k + \frac{(\underline{p}_k - L_k \underline{u}_k) \underline{u}_k^T}{\underline{u}_k^T \underline{u}_k} \right] \left[ L_k^T + \frac{\underline{u}_k (\underline{p}_k - L_k \underline{u}_k)^T}{\underline{u}_k^T \underline{u}_k} \right]$$





# Recovering DFP Update

$$\begin{aligned}
 H_{k+1} &= H_k + \frac{L_k \underline{u}_k (\underline{p}_k - L_k \underline{u}_k)^\top}{\underline{u}_k^\top \underline{u}_k} + \frac{(\underline{p}_k - L_k \underline{u}_k) \underline{u}_k^\top L_k^\top}{\underline{u}_k^\top \underline{u}_k} + \frac{(\underline{p}_k - L_k \underline{u}_k) (\underline{p}_k - L_k \underline{u}_k)^\top}{\underline{u}_k^\top \underline{u}_k} \\
 &= H_k + \frac{\underline{p}_k \underline{p}_k^\top}{\underline{p}_k^\top \underline{q}_k} - \frac{L_k \underline{u}_k \underline{u}_k^\top L_k^\top}{\underline{u}_k^\top \underline{u}_k}
 \end{aligned}$$

$$\begin{aligned}
 \underline{u}_k^\top \underline{u}_k &= \underline{p}_k^\top \underline{q}_k \\
 \underline{u}_k &= \Theta_k L_k^\top \underline{q}_k \\
 H_k &= L_k L_k^\top
 \end{aligned}$$

$$H_{k+1} = H_k + \frac{\underline{p}_k \underline{p}_k^\top}{\underline{p}_k^\top \underline{q}_k} - \frac{H_k \underline{q}_k \underline{q}_k^\top H_k}{\underline{q}_k^\top H_k \underline{q}_k}$$

Davidon-Fletcher-Powell (DFP) method (or) Inverse Positive Definite Secant Update

- All we need to do to ensure a descent direction ( $H_{k+1} > 0$ ) is

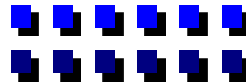
$$\underline{p}_k^\top \underline{q}_k > 0 \Rightarrow \underline{g}_{k+1}^\top \underline{d}_k > \underline{g}_k^\top \underline{d}_k$$

Line search can be inexact.

- Because of round-off errors,  $H_{k+1}$  can still go indefinite. Can rewrite  $H_{k+1}$  as:

$$H_{k+1} = \left[ I - \frac{H_k \underline{q}_k \underline{q}_k^\top}{\underline{q}_k^\top H_k \underline{q}_k} \right] H_k \left[ I - \frac{H_k \underline{q}_k \underline{q}_k^\top}{\underline{q}_k^\top H_k \underline{q}_k} \right]^\top + \frac{\underline{p}_k \underline{p}_k^\top}{\underline{p}_k^\top \underline{q}_k} + \frac{H_k \underline{q}_k \underline{q}_k^\top H_k}{\underline{q}_k^\top H_k \underline{q}_k}$$

Similar to Joseph's form of covariance equation





# Recovering BFGS Update

- We can obtain Hessian approximation,  $B_k$  via:

$$L_k L_k^T = H_k \iff B_k = \hat{L}_k \hat{L}_k^T \quad \underline{q}_k \leftrightarrow \underline{p}_k \quad \underline{p}_k \leftrightarrow \underline{q}_k \quad L_k, S_k \leftrightarrow \hat{L}_k, \hat{S}_k$$

$$\hat{S}_{k+1} = \hat{L}_k + \frac{(\underline{q}_k - \hat{L}_k \hat{u}_k) \hat{u}_k^T}{\hat{u}_k^T \hat{u}_k}, \quad \hat{u}_k = \sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{p}_k^T B_k \underline{p}_k}} \cdot \hat{L}_k^T \underline{p}_k$$

$$B_{k+1} = B_k + \frac{\underline{q}_k \underline{q}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{B_k \underline{p}_k \underline{p}_k^T B_k}{\underline{p}_k^T B_k \underline{p}_k} \quad (*)$$

Broyden-Fletcher-Goldfarb-Shanno (BFGS) formula (or) positive definite secant update

- Found to be the best update in practice.**  $\underline{d}_k$  from  $\hat{L}_k \hat{L}_k^T \underline{d}_k = -\underline{g}_k$
- We arbitrarily modify  $H_{k+1}$  update as:

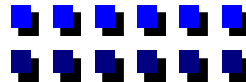
$$H_{k+1} = B_{k+1}^{-1} \text{ of } (*) \text{ modified}$$

$$= H_k + \frac{\underline{p}_k \underline{p}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{H_k \underline{q}_k \underline{q}_k^T H_k}{\underline{q}_k^T H_k \underline{q}_k} + \phi \left( \underline{q}_k^T H_k \underline{q}_k \right) \cdot \underline{r}_k \underline{r}_k^T, \quad \underline{r}_k = \frac{\underline{p}_k}{\underline{p}_k^T \underline{q}_k} - \frac{1}{\underline{q}_k^T H_k \underline{q}_k} \cdot H_k \underline{q}_k$$

$\phi = 1$  actual  $B_{k+1}^{-1}$  BFGS inverse PD secant update

$\phi = 0$  DFP

Note that  $\underline{r}_k^T \underline{q}_k = 0 \Rightarrow H_{k+1} \underline{q}_k = \underline{p}_k$  is satisfied.





# Analogy with Control & Estimation

- Similarly, we can arbitrarily modify  $B_{k+1}$  update as:

$B_{k+1} = H_{k+1}^{-1}$  of DFP modified

$$= B_k + \frac{\underline{q}_k \underline{q}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{B_k \underline{p}_k \underline{p}_k^T B_k}{\underline{p}_k^T B_k \underline{p}_k} + \phi \left( \underline{p}_k^T B_k \underline{p}_k \right) \underline{w}_k \underline{w}_k^T$$

$$\underline{w}_k = \frac{\underline{q}_k}{\underline{p}_k^T \underline{q}_k} - \frac{B_k \underline{p}_k}{\underline{p}_k^T B_k \underline{p}_k}, \quad \underline{w}_k^T \underline{p}_k = 0 \Rightarrow B_{k+1} \underline{p}_k = \underline{q}_k \quad \phi = \begin{cases} 1 & H_k^{-1} \text{ update of DFP} \\ 0 & \text{BFGS (best)} \end{cases}$$

- Expressions for  $B_{k+1}$  and  $H_{k+1}$  are similar to those that arise in filtering (with perfect observations) and linear quadratic control problems (with no control weightings).

## Filtering

$$\underline{x}_{k+1} = \underline{x}_k + \underline{p}_k \underline{w}_k; \quad \underline{w}_k = N \left( 0, \frac{1}{\underline{p}_k^T \underline{q}_k} \right)$$

$$\underline{y}_k = \underline{q}_k^T \underline{x}_k \quad (\text{NO M/S NOISE})$$

$$\hat{\underline{x}}_{k+1} = \hat{\underline{x}}_k + K_k \left[ y_k - \underline{q}_k^T \hat{\underline{x}}_k \right]$$

$$K_k = \frac{\underline{\Sigma}_k \underline{q}_k}{\underline{q}_k^T \underline{\Sigma}_k \underline{q}_k}, \quad \underline{\Sigma}_{k+1} = \underline{\Sigma}_k + \frac{\underline{p}_k \underline{p}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{\underline{\Sigma}_k \underline{q}_k \underline{q}_k^T \underline{\Sigma}_k}{\underline{q}_k^T \underline{\Sigma}_k \underline{q}_k}$$

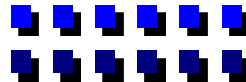
## Control

$$\underline{x}_{k+1} = \underline{x}_k + \underline{q}_k \underline{u}_k$$

$$y_k = \underline{p}_k^T \underline{x}_k$$

$$\min \left\{ \sum_{k=0}^N \frac{1}{\underline{p}_k^T \underline{q}_k} y_k^2 \right\}$$

$$\underline{P}_{N-k} = \underline{P}_{N-k+1} + \frac{\underline{p}_{N-k+1} \underline{p}_{N-k+1}^T}{\underline{p}_{N-k+1}^T \underline{q}_{N-k+1}} - \frac{\underline{P}_{N-k+1} \underline{q}_{N-k+1} \underline{q}_{N-k+1}^T \underline{P}_{N-k+1}}{\underline{q}_{N-k+1}^T \underline{P}_{N-k+1} \underline{q}_{N-k+1}}$$





# How to Traingularize Matrix S?

□ Problem 2: How to get  $L_{k+1}$  from  $S_{k+1} \ni S_{k+1} = L_{k+1} T_{k+1}^T$  &  $T_{k+1}$  orthogonal?

- $\exists$  methods for transforming a matrix  $A$  into a product of a lower triangular matrix and an orthogonal matrix.

- Householder, Givens, etc. (ECE 6435). Complexity  $O(n^3)$ .

- But,  $S_{k+1}$  is a rank one perturbation of a lower  $\Delta$  matrix,  $L_k$ , i.e.,

$$S_{k+1} = L_k + \frac{(\underline{p}_k - L_k \underline{u}_k) \underline{u}_k^T}{\underline{u}_k^T \underline{u}_k}$$

- For this problem, we can devise an  $O(n^2)$  algorithm based on

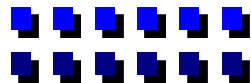
GIVENS ROTATIONS

□ What we have is:  $L_k^T = I \cdot L_k^T \Rightarrow T_k = I, L_k^T = \text{upper } \Delta$ .

We want  $S_{k+1}^T = T_{k+1}^T L_{k+1}^T \rightarrow$  QR decomposition of  $S_{k+1}^T$

- Don't actually need to store  $T_{k+1}$  since  $H_{k+1} = S_{k+1} S_{k+1}^T = L_{k+1} L_{k+1}^T$ .
- The generic problem we want to solve is to find the QR decomposition of

$$S^T = L^T + \underline{u} \underline{y}^T \quad \underline{y} \sim \frac{\underline{p}_k - L_k \underline{u}_k}{\underline{u}_k^T \underline{u}_k}$$





# Triangularization via Givens

- The algorithm involves two steps:

Step 1: Find an orthogonal transform  $T_1$  such that  $U_H = T_1^T [L^T + \underline{u}\underline{y}^T]$ .

$$\begin{bmatrix} \diagup & & & & \\ & \diagdown & & & \\ & & \diagdown & & \\ & & & \diagdown & \\ & & & & \diagdown \end{bmatrix} + \begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix} \underline{T}_1 \begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ 0 & x & x & x & x \\ 0 & 0 & x & x & x \\ 0 & 0 & 0 & x & x \end{bmatrix}$$

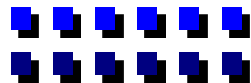
$L^T$                        $\underline{u}\underline{y}^T$                        $U_H$  upper Hessenberg matrix

Step 2: Find a second orthogonal transformation

$$T_2 \ni \tilde{L}^T = T_2^T U_H = T_2^T T_1^T [L^T + \underline{u}\underline{y}^T].$$

$$\begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ 0 & x & x & x & x \\ 0 & 0 & x & x & x \\ 0 & 0 & 0 & x & x \end{bmatrix} \underline{T}_2 \begin{bmatrix} x & x & x & x & x \\ 0 & x & x & x & x \\ 0 & 0 & x & x & x \\ 0 & 0 & 0 & x & x \\ 0 & 0 & 0 & 0 & x \end{bmatrix}$$

$$T = T_1 T_2 \quad \& \quad \tilde{L}^T = T^T [L^T + \underline{u}\underline{y}^T]$$







# Upper Hessenberg via Givens -1

## Step 1:

1.1 Find  $T_1^T \ni T_1^T \underline{u} = [\alpha \ 0 \ \dots \ 0]^T = \alpha \underline{e}_1$ ;  $\alpha = \pm\sqrt{\underline{u}^T \underline{u}}$

1.2 Apply  $T_1^T$  to  $L^T + \underline{u}\underline{y}^T \Rightarrow \underbrace{T_1^T L^T}_{\text{upper Hessenberg}} + \underbrace{\alpha \underline{e}_1 \underline{y}^T}_{\bullet \curvearrowright} = U_H.$

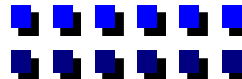
$$\left[ \begin{array}{ccc} \leftarrow & \alpha \underline{y}^T & \rightarrow \\ \leftarrow & 0 & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & 0 & \rightarrow \end{array} \right]$$

– The tool to accomplish step 1 is GIVENS rotations.

$$J(i,k) \triangleq J(i,k,\theta) = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & & & & \\ & & & c & s & \\ & & & -s & c & \\ & & & & & \ddots & \\ & & & & & & 1 \end{bmatrix} \quad \left. \begin{array}{l} c = \cos \theta \\ s = \sin \theta \end{array} \right\} \Rightarrow \sqrt{c^2 + s^2} = 1$$

$$J^{-1} = J^T = J(i,k,-\theta)$$

$$y = J\underline{x} \Rightarrow y_i = cx_i + sx_k; \quad y_k = -sx_i + cx_k; \quad y_j = x_j \quad \forall j \neq i,k$$





# Upper Hessenberg via Givens -2

- If you want to zero out  $y_k$ , select  $c = \frac{x_i}{\sqrt{x_i^2 + x_k^2}}$ ,  $s = \frac{x_k}{\sqrt{x_i^2 + x_k^2}}$

– So, to transform  $\underline{u} \rightarrow \alpha \underline{e}_1$  :

$$\bullet \overbrace{J(1,2) \dots J(n-2,n)^{-1} J(n-1,n)}^{T_1^T} \underline{u} = \alpha \underline{e}_1 = \begin{bmatrix} u_1 & 0 & \dots & 0 \end{bmatrix}^T = u_1 \underline{e}_1$$

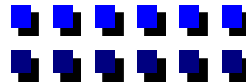
- Whatever you do to  $\underline{u}$ , do the same to  $L^T$ , i.e.,

$$J(1,2) \dots J(n-2,n) J(n-1,n) L^T = T_1^T L^T = \text{upper Hessenberg}$$

- Actually can do operations directly on  $L$

$$\Rightarrow \text{lower Hessenberg} \Rightarrow L_{H_1} = L T_1 = L J^T(n-1,n) \dots J^T(1,2)$$

$$\begin{bmatrix} \times \\ \times \\ \times \\ \times \end{bmatrix} \xrightarrow{J(3,4)} \begin{bmatrix} \times \\ \times \\ \otimes \\ 0 \end{bmatrix} \xrightarrow{J(2,3)} \begin{bmatrix} \times \\ \otimes \\ 0 \\ 0 \end{bmatrix} \xrightarrow{J(1,2)} \begin{bmatrix} \otimes \\ 0 \\ 0 \\ 0 \end{bmatrix}$$





# Step 1 Algorithm via Givens

## Algorithm:

FOR  $i = n-1, n-2, \dots, 1$  DO

$$c = \frac{u_i}{\sqrt{u_i^2 + u_{i+1}^2}}$$

$$s = \frac{u_{i+1}}{\sqrt{u_i^2 + u_{i+1}^2}}$$

$$u_i \leftarrow cu_i + su_{i+1}$$

$$u_{i+1} \leftarrow 0$$

FOR  $j = n, n-1, \dots, i$  DO

$$a = cl_{ji} + sl_{j,i+1}$$

$$b = -sl_{ji} + cl_{j,i+1}$$

$$l_{ji} \leftarrow a$$

$$l_{j,i+1} \leftarrow b$$

END

END

FOR  $i = 1, n$  DO

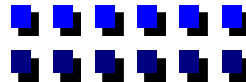
$$l_{i1} \leftarrow l_{i1} + u_1 y_i$$

END

only 2  
columns of  $L$   
( $j, i+1$ ) are  
affected.

$$L_{H_1} = LT_1$$

$$L_{H_1} + \alpha \underline{y}_1 e^T = L_H = U_H^T$$





# Step 2 Algorithm via Givens

$$\begin{bmatrix} x & 0 & 0 & 0 \\ x & x & 0 & 0 \\ x & x & x & 0 \\ x & x & x & x \end{bmatrix} \xrightarrow{J^T(3,4)} \begin{bmatrix} x & 0 & 0 & 0 \\ x & x & 0 & 0 \\ x & x & \Delta & \Delta \\ x & x & \Delta & \Delta \end{bmatrix} \xrightarrow{J^T(2,3)} \begin{bmatrix} x & 0 & 0 & 0 \\ x & \square & \square & 0 \\ x & \square & \square & \Delta \\ x & \square & \square & \Delta \end{bmatrix} \xrightarrow{J^T(2,1)} \begin{bmatrix} \diamond & \diamond & 0 & 0 \\ \diamond & \diamond & \square & 0 \\ \diamond & \diamond & \square & \Delta \\ \diamond & \diamond & \square & \Delta \end{bmatrix}$$

- Step 2: Can we use Givens rotations again to transform Hessenberg to desired lower triangular form? **YES**

FOR  $j = n, n-1, \dots, 2$  DO

$$c = \frac{l_{jj}}{\sqrt{l_{j-1,j}^2 + l_{jj}^2}}$$

$$s = \frac{-l_{j-1,j}}{\sqrt{l_{jj}^2 + l_{j-1,j}^2}}$$

$$l_{jj} \leftarrow \sqrt{l_{jj}^2 + l_{j-1,j}^2}$$

$$l_{j-1,j} \leftarrow 0$$

FOR  $i = 1, 2, \dots, j-1$  DO

$$a \leftarrow cl_{j-1,i} + sl_{ji}$$

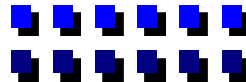
$$l_{ji} \leftarrow -sl_{j-1,i} + cl_{ji}$$

$$l_{j-1,i} \leftarrow a$$

} rows  $j-1$  &  $j$  are affected

END

END





# References

## □ References:

1. D. Goldfarb, “Factorized variable metric methods for unconstrained minimization,” Mathematics of Computation, vol. 30, NO. 136, Oct. 1976, pp. 796-811
2. B. Fletcher and M. J. D. Powell, “On the modifications of  $LDL^T$  factorizations,” Mathematics of Computation, vol. 28, NO. 128, Oct. 1976, pp. 1067-1087.
3. P. E. Gill, G. H. Golub, W. Murray, and M. A. Saunders, “Methods for modifying matrix factorizations,” Mathematics of Computation, vol. 28, NO. 126, April 1974, pp. 505-535.



# Quadratic Termination of QN Methods

## 1. Quadratic Termination: DFP method yields

- $Q$ -conjugate  $\underline{p}_i$  ( $\Rightarrow Q$ -conjugate  $\underline{d}_i$  since  $\underline{p}_i = \alpha_i \underline{d}_i$ )
- $H_{k+1} Q \underline{p}_i = \underline{p}_i, \quad 0 \leq i \leq k$
- $H_n = Q^{-1}$ 
  - In the quadratic case:

$$\underline{q}_k = \underline{g}_{k+1} - \underline{g}_k = Q(\underline{x}_{k+1} - \underline{x}_k) = Q \underline{p}_k$$

- Also, quasi-Newton condition yields:

$H_{k+1} \underline{q}_k = \underline{p}_k \Rightarrow H_{k+1} Q \underline{p}_k \Rightarrow \underline{p}_k$  is an eigenvector of  $H_{k+1} Q$  with a unity eigenvalue

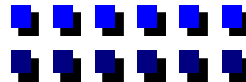
- We prove  $Q$ -conjugacy of  $\underline{p}_i$  via induction. Assume that  $\underline{p}_0, \underline{p}_1, \dots, \underline{p}_k$  are  $Q$ -conjugate. We prove that  $\underline{p}_0, \underline{p}_1, \dots, \underline{p}_{k+1}$  are  $Q$ -conjugate. We know that

$$\underline{g}_{k+1} = \underline{g}_{i+1} + Q[\underline{p}_{i+1} + \underline{p}_{i+2} + \dots + \underline{p}_k]$$

$$\underline{p}_i^T \underline{g}_{k+1} = \underline{p}_i^T \underline{g}_{i+1} = 0 \quad \forall 0 \leq i \leq k \quad \because \text{of } Q\text{-conjugacy}$$

$$\Rightarrow \underline{p}_i^T \underline{g}_{k+1} = 0 \Rightarrow \underline{p}_i^T \underbrace{Q H_{k+1} \underline{g}_{k+1}}_{\propto \underline{p}_{k+1}} = 0 \quad \because H_{k+1} Q \underline{p}_i = \underline{p}_i$$

$$\text{Since } \underline{p}_{k+1} = -\alpha_k H_{k+1} \underline{g}_{k+1} \Rightarrow \underline{p}_i^T Q \underline{p}_{k+1} = 0 \quad 0 \leq i \leq k$$





# Memoryless Quasi-Newton Method

- Also, for  $k$

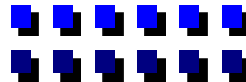
$$\underline{q}_{k+1}^T H_{k+1} Q \underline{p}_i = \underline{q}_{k+1}^T \underline{p}_i = \underline{p}_{k+1}^T Q \underline{p}_i = 0 \quad \forall 0 \leq i \leq k$$

$$\Rightarrow H_{k+2} Q \underline{p}_i = H_{k+1} Q \underline{p}_i + \frac{\overbrace{\underline{p}_{k+1} \underline{p}_{k+1}^T Q \underline{p}_i}^{=0}}{\underline{p}_{k+1}^T \underline{q}_{k+1}} - \frac{H_{k+1} \underline{q}_{k+1} \left( \overbrace{\underline{q}_{k+1}^T H_{k+1} Q \underline{p}_i}^{=0} \right)}{\underline{q}_{k+1}^T H_{k+1} \underline{q}_{k+1}} = \underline{p}_i \quad \forall 0 \leq i \leq k$$

- Since  $\underline{p}_k$ 's are  $Q$ -orthogonal, DFP is a conjugate gradient method.
- $H_0 = I \Rightarrow$  DFP  $\Rightarrow$  conjugate gradient method
- $H_k = I \quad \forall k \Rightarrow$ 
  - BFGS inverse positive definite secant update is  $\equiv$  conjugate gradient method.
  - Also known as memoryless Quasi-Newton.
  - Recall:

$$H_{k+1} = H_k + \frac{\underline{p}_k \underline{p}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{H_k \underline{q}_k \underline{q}_k^T H_k}{\underline{q}_k^T H_k \underline{q}_k} + \underline{q}_k^T H_k \underline{q}_k \underline{r}_k \underline{r}_k^T, \quad \underline{r}_k = \frac{\underline{p}_k}{\underline{p}_k^T \underline{q}_k} - \frac{H_k \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}$$

$$\Rightarrow H_{k+1} = I + \frac{\underline{p}_k \underline{p}_k^T}{\underline{p}_k^T \underline{q}_k} - \frac{\underline{q}_k \underline{q}_k^T}{\underline{q}_k^T \underline{q}_k} + \underline{q}_k^T \underline{q}_k \underline{r}_k \underline{r}_k^T = I - \frac{\underline{q}_k \underline{p}_k^T + \underline{p}_k \underline{q}_k^T}{\underline{p}_k^T \underline{q}_k} + \left[ 1 + \frac{\underline{q}_k^T \underline{q}_k}{\underline{p}_k^T \underline{q}_k} \right] \frac{\underline{p}_k \underline{p}_k^T}{\underline{p}_k^T \underline{q}_k}$$





# Memoryless QN & Sorensen-Wolfe CG

$$\begin{aligned} \underline{d}_{k+1} &= -H_{k+1} \underline{g}_{k+1} \\ &= -\underline{g}_{k+1} + \frac{\underline{q}_k \underline{p}_k^T \underline{g}_{k+1} + \underline{p}_k \underline{q}_k^T \underline{g}_{k+1}}{\underline{p}_k^T \underline{q}_k} + \left[ 1 + \frac{\underline{q}_k \underline{q}_k^T}{\underline{p}_k^T \underline{q}_k} \right] \underline{p}_k \underline{p}_k^T \underline{g}_{k+1} \end{aligned}$$

- For exact line searches,  $\underline{p}_k^T \underline{g}_{k+1} = 0$

$$\Rightarrow \boxed{\underline{d}_{k+1} = -\underline{g}_{k+1} + \frac{\underline{q}_k^T \underline{g}_{k+1}}{\underline{d}_k^T \underline{q}_k} \underline{d}_k} \longleftarrow \text{Sorensen-Wolfe formula}$$

## □ Convergence Properties: Reset $H_n = I$ , i.e., every $n$ steps.

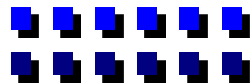
- Since approximate  $\nabla^2 f(\underline{x}_k)$  near minimum, get superlinear convergence.

$$\Rightarrow \beta = \sup \lim_{k \rightarrow \infty} \frac{\|\underline{x}_{k+1} - \underline{x}'\|}{\|\underline{x}_k - \underline{x}'\|} = 0$$

- In general, BFGS positive definite secant update is better than DFP.
- Can improve convergence via scaling.

## □ Scaling: Convergence depends on $\lambda_i(H_k^{1/2} Q H_k^{1/2}) = \lambda_i(H_k Q)$

- Recall that for quadratic functions  $H_{k+1} Q \underline{p}_i = \underline{p}_i \quad \forall 0 \leq i \leq k$







# Scaling in Quasi-Newton Methods - 1

## Scaling cont.,

- At each step of DFP (or Broyden), we move one eigen value at a time to unity.
- $\kappa(H_k Q)$  may be worse than  $\kappa(Q)$ .

**Idea:** Scale  $H_k$  by a scalar  $\gamma_k$  so that the eigen values of  $H_k Q$  are spread both below and above unity.

- If  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  are eigenvalues of  $H_k Q$ , we want to multiply  $H_k$  by  $\gamma_k$  where  $\lambda_1 \leq \frac{1}{\gamma_k} \leq \lambda_n$ .

$$\Rightarrow \lambda_i [\gamma_k H_k Q] \Rightarrow \gamma_k \lambda_1 \leq 1 \leq \gamma_k \lambda_n$$

Recall:

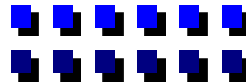
$$\frac{\underline{q}_k^T H_k \underline{q}_k}{\underline{p}_k^T \underline{q}_k} = \frac{\underline{p}_k^T Q H_k Q \underline{p}_k}{\underline{p}_k^T Q \underline{p}_k} = \frac{\underline{r}_k^T Q^{1/2} H_k Q^{1/2} \underline{r}_k}{\underline{r}_k^T \underline{r}_k} \quad \text{where } \underline{r}_k = Q^{1/2} \underline{p}_k$$

- Since  $\lambda_i(Q^{1/2} H_k Q^{1/2}) = \lambda_i(H_k^{1/2} Q H_k^{1/2}) = \lambda_i(H_k Q)$

$$\lambda_i(AB) = \lambda_i(BA)$$

$$\lambda_1(H_k Q) \leq \frac{\underline{q}_k^T H_k \underline{q}_k}{\underline{p}_k^T \underline{q}_k} \leq \lambda_n(H_k Q)$$

$$\text{So, use } \gamma_k = \frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}$$





## Scaling in Quasi-Newton Methods - 2

□ The scaled quasi-Newton method then works as follows:

- \*  $k = 0, H_0 = I$
- †  $\underline{d}_k = -H_k \underline{g}_k$
- Perform line search to find  $\alpha_k = \arg \min_{\alpha} f(\underline{x}_k + \alpha \underline{d}_k)$
- If  $k = n$  Go back to \*

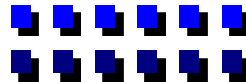
$$\text{else } \tilde{H}_k = H_k \frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T H_k \underline{q}_k}$$

update  $H_k$  and go back to †.

- For square root implementation  $\tilde{L}_k = \sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{q}_k^T L_k L_k^T \underline{q}_k}} \cdot L_k = \text{scaled } L_k$

- For BFGS positive definite secant update, the scaling factor is  $\sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{p}_k^T B_k \underline{p}_k}}$ .

$$\text{Then } \tilde{B}_k = B_k \cdot \frac{\underline{q}_k^T \underline{p}_k}{\underline{p}_k^T B_k \underline{p}_k} \quad \& \quad \tilde{L}_k = \sqrt{\frac{\underline{p}_k^T \underline{q}_k}{\underline{p}_k^T \hat{L}_k \hat{L}_k^T \underline{p}_k}} \cdot \hat{L}_k$$





# Neural Networks: Least Square Problems

## □ Least Square Problems

$$\min f(\underline{x}) = \frac{1}{2} \sum_{i=1}^m g_i^2(\underline{x})$$

$$\nabla_{\underline{x}} f(\underline{x}) = \sum_{i=1}^m \nabla_{\underline{x}} g_i(\underline{x}) g_i(\underline{x}) = \nabla_{\underline{x}} \underline{g}(\underline{x}) \underline{g}(\underline{x}) = J^T(\underline{x}) \underline{g}(\underline{x})$$

$$\nabla^2 f(\underline{x}) = \nabla_{\underline{x}} \underline{g}(\underline{x}) \nabla_{\underline{x}} \underline{g}^T(\underline{x}) + \sum_{i=1}^m \nabla^2 g_i(\underline{x}) g_i(\underline{x})$$

## □ Examples:

$$1) f(\underline{x}) = \frac{1}{2} \sum_{i=1}^m (b_i - \underline{a}_i^T \underline{x})^2 = \frac{1}{2} \sum_{i=1}^m e_i^2$$

$$\nabla_{\underline{x}} f(\underline{x}) = - \sum_{i=1}^m (b_i - \underline{a}_i^T \underline{x}) \underline{a}_i = - \sum_{i=1}^m e_i \underline{a}_i$$

$$\nabla_{\underline{x}} e_i = \underline{a}_i$$

$$\nabla^2 f(\underline{x}) = \sum_{i=1}^m \underline{a}_i \underline{a}_i^T$$



# Examples of Least Squares Problems - 1

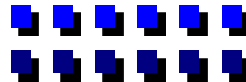
## Examples (continued) :

$$2) f(\underline{x}) = \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=1}^n e^{t_i x_j} - y_i \right)^2; \quad \{t_i, y_i\}_{i=1}^m$$

$$\nabla_{\underline{x}} f(\underline{x}) = \underbrace{\begin{bmatrix} t_1 e^{t_1 x_1} & t_2 e^{t_2 x_1} & t_m e^{t_m x_1} \\ t_1 e^{t_1 x_2} & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ t_1 e^{t_1 x_n} & t_2 e^{t_2 x_n} & t_m e^{t_m x_n} \end{bmatrix}}_{\nabla_{\underline{g}}} \begin{bmatrix} \sum_{i=1}^n e^{t_1 x_j} - y_1 \\ \sum_{i=1}^n e^{t_2 x_j} - y_2 \\ \vdots \\ \sum_{i=1}^n e^{t_m x_j} - y_m \end{bmatrix}$$

$$\nabla^2 f(\underline{x}) = \nabla_{\underline{g}} \nabla_{\underline{g}}^T + \sum_{i=1}^m g_i(\underline{x}) \nabla^2 g_i(\underline{x})$$

$$\nabla^2 g_i(\underline{x}) = \begin{bmatrix} t_i^2 e^{t_i x_1} & & & \\ & t_i^2 e^{t_i x_2} & & \\ & & \ddots & \\ & & & t_i^2 e^{t_i x_n} \end{bmatrix}$$

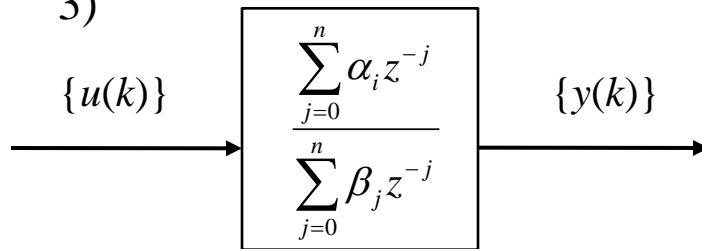




# Examples of Least Squares Problems - 2

## Examples (continued):

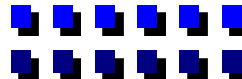
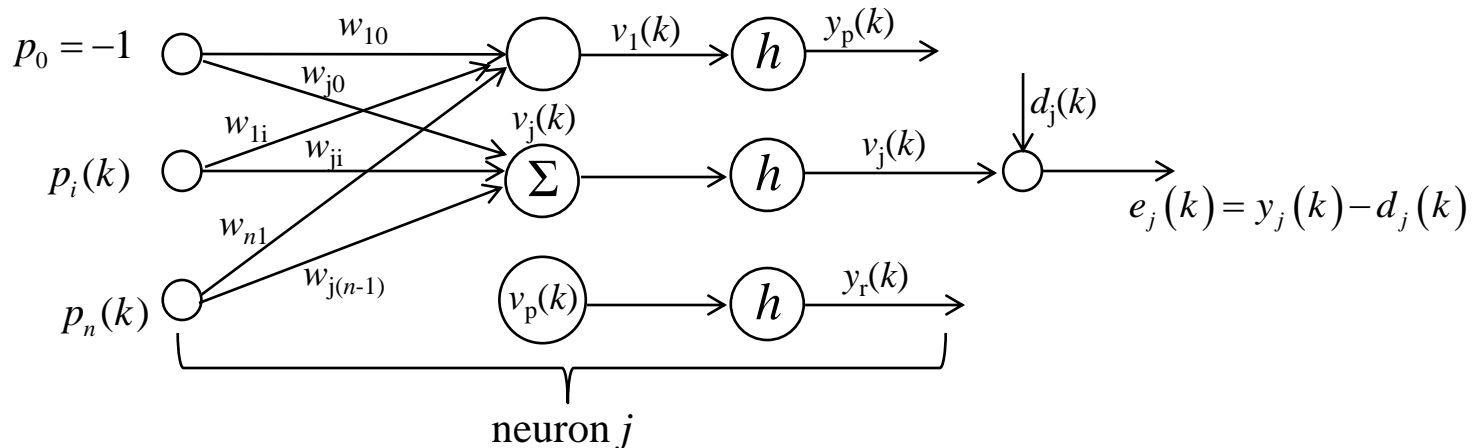
3)



We want to estimate  $\{\alpha_i\}$  and  $\{\beta_i\}$  from  $\{u(k), y(k)\}_{k=1}^m$   $m \gg 2n$

$$f(\underline{\alpha}, \underline{\beta}) = \frac{1}{2} \sum_{k=n}^m \left( \sum_{j=0}^n \alpha_j u_{k-j} - \sum_{j=0}^n \beta_j y_{k-j} \right)^2$$

4) Single layer neural networks: Given pairs of  $\{\underline{p}(k), \underline{d}(k)\}_{k=1}^m$ , find network weight  $W$  to minimize  $\frac{1}{2m} \sum_{i=1}^r e_i^2(k) \rightarrow$  average error.





# Examples of Least Squares Problems - 3

## □ Examples (continued):

$$4) \text{ cont... } v_j(k) = \sum_{i=0}^n w_{ji} p_i(k) = \underline{w}_j^T \underline{p}(k)$$

$$y_j(k) = h(v_j(k))$$

$$h(v_j) = \begin{cases} \frac{1}{1 + e^{-\sigma v_j}} & \in (0,1) \text{ logistic function} \\ \frac{2}{\pi} \tan^{-1}(v_j) & \in (-1,1) \\ \frac{1 - e^{-\sigma v_j}}{1 + e^{-\sigma v_j}} & \text{bipolar for } \sigma = 1, \text{ tanh for } \sigma = 2 \end{cases}$$

$$f(W) = \frac{1}{2m} \sum_{k=1}^m \sum_{j=1}^r (d_j(k) - h(\underline{w}_j^T \underline{p}(k)))^2 = \frac{1}{2} \sum_{j=1}^r f_j(\underline{w}_j)$$

- Since the cost function is separable, we can minimize this function for each output separately.

$$f_j(\underline{w}_j) = \frac{1}{2m} \sum_{k=1}^m (d_j(k) - h(\underline{w}_j^T \underline{p}(k)))^2$$

- So, the canonical problem is:  $f(\underline{w}) = \frac{1}{2} \sum_{k=1}^m (d(k) - h(\underline{w}_j^T \underline{p}(k)))^2$



# General Neural Networks

- The Adaline (Adaptive Linear Neuron) Widrow-Hoff (1962)  
LMS algorithm (Least Mean Square)

$$h(\underline{w}^T \underline{p}) = \underline{w}^T \underline{p}$$

$$f(\underline{w}) = \frac{1}{2m} \sum_{k=1}^m \underbrace{(d(k) - \underline{w}^T \underline{p}(k))^2}_{e^2(k)}$$

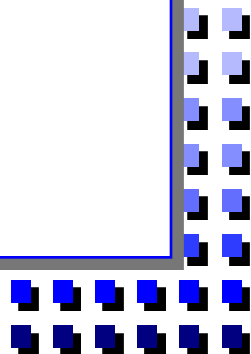
$$\nabla f(\underline{w}) = -\frac{1}{m} \sum_{k=1}^m (d(k) - \underline{w}^T \underline{p}(k)) \underline{p}(k) = -\frac{1}{m} \sum_{i=1}^m e(k) \underline{p}(k)$$

$$\nabla^2 f(\underline{w}) = \frac{1}{m} \sum_{k=1}^m \underline{p}(k) \underline{p}^T(k)$$

- General Neural Networks

- Can approximate any I/O map accurately
- Adapt to stochastic environments
- Fault-tolerance & VLSI implementability
- Theoretically, one hidden layer is sufficient

$$h(x) = \frac{1}{1 + e^{-\sigma x}}, \quad h'(x) = \sigma h(x)(1 - h(x))$$





# Multilayer Perceptrons

$$\begin{aligned} \nabla_{\underline{w}} f(\underline{w}) &= -\sum_{k=1}^m e(k) h'(\underline{w}^T \underline{p}(k)) \underline{p}(k) \\ &= -\sigma \sum_{k=1}^m \left[ d(k) - h(\underline{w}^T \underline{p}(k)) \right] h(\underline{w}^T \underline{p}(k)) \left[ 1 - h(\underline{w}^T \underline{p}(k)) \right] \underline{p}(k) \end{aligned}$$

Note that:

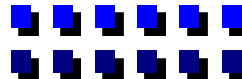
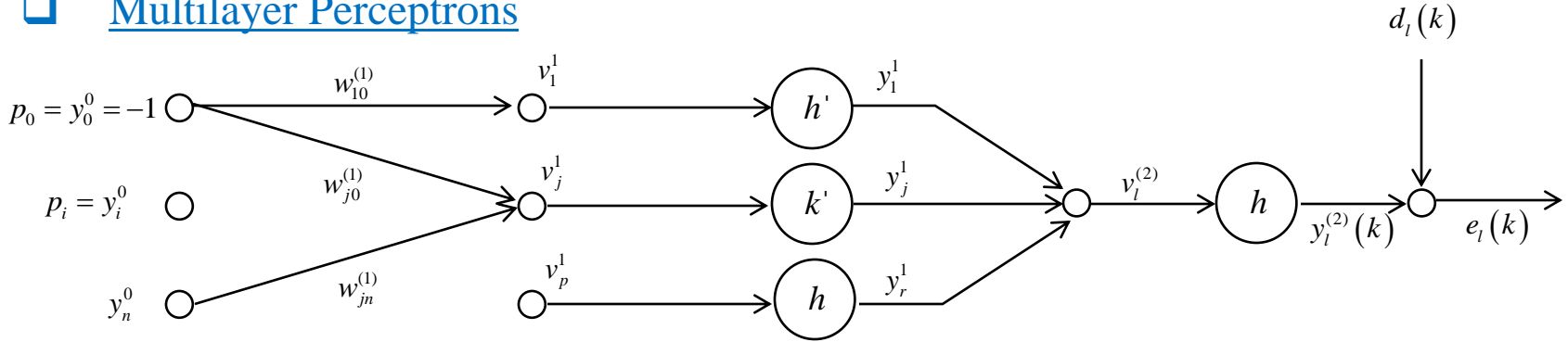
$$\frac{\partial f_k}{\partial w_i} = \underbrace{-e(k) h'(v(k))}_{\delta(k)} p_i(k)$$

Local computation  
local gradient

local gradient of neuron  $j$

- In general,  $\frac{\partial f}{\partial w_{ji}(k)} = -\delta_j(k) p_i(k)$ . Second derivatives are a little more complex (see B.D. Ripley, *Pattern recognition & Neural Networks*, 1996 or Bishop, 2006)

## □ Multilayer Perceptrons







# Gradient Computation

$$y_j^1 = h^1 \left[ \left( \underline{w}_j^1 \right)^T \underline{y}^0 \right]; \quad y_k^2 = h^2 \left[ \left( \underline{w}_k^2 \right)^T \underline{y}^1 \right] = h^2 \left( \begin{array}{c} \underline{w}_k^T \\ \vdots \\ h \left[ \left( \underline{w}_n^1 \right)^T \underline{y}^0 \right] \end{array} \right)$$

**Key:** Can evaluate gradient recursively backwards

$$\frac{\partial f}{\partial w_{ji}^q(k)} = \delta_j^q(k) y_i^{(q-1)}(k); \quad \delta_j^q(k) = h^1(v_j^q(k)) \sum_l \delta_l^{q+1}(k) w_{lj}^{q+1}(k)$$

$$\text{end: } \delta_j^L(k) = -e(k) h^L(v^L(k)) \quad (\text{Last Layer})$$

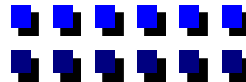
- Can use SD, CG, QN, Newton, Modified Newton, etc.
- Gauss-Newton neglects second order terms.

$$\nabla^2 f = \left[ \nabla \underline{g} \nabla \underline{g}^T + \sum_{i=1}^m \nabla^2 g_i g_i \right] \approx \nabla \underline{g} \nabla \underline{g}^T$$

$$\underline{x}_{k+1} = \underline{x}_k - \left[ \nabla \underline{g}(\underline{x}_k) \nabla \underline{g}^T(\underline{x}_k) \right]^{-1} \nabla \underline{g}(\underline{x}_k) \underline{g}(\underline{x}_k)$$

– If not invertible:

$$\nabla \underline{g}(\underline{x}_k) \nabla \underline{g}^T(\underline{x}_k) + \Delta_k I \quad \text{Levenberg-Marquardt}$$





# Incremental Gradient Algorithm -1

## □ Some neat implementations:

- Recursive  $k = 1, 2, \dots, m$  (Incremental Gradient Methods)

$$f(\underline{x}) = \sum_{l=1}^m f_l(\underline{x}) = \sum_{l=1}^m g_l^2(\underline{x}) \quad g_1, g_2, \dots, g_m \text{ come in sequentially}$$

$$\underline{x}_{k+1} = \underline{x}_k - \alpha_k \sum_{l=1}^m \nabla g_l(\underline{x}_k) g_l(\underline{x}_k); \quad \Psi_0 = \underline{x}_k$$

FOR  $l = 1, m$  DO

$$\Psi_l = \Psi_{l-1} - \alpha_k^{(l)} \nabla g_l(\Psi_{l-1}) g_l(\Psi_{l-1})$$

END

- This is precisely what happens in LMS, Recursive Least Squares and Neural Networks

- Use constant step size  $\alpha_k^{(l)} = \alpha_k = \min\left(\gamma, \frac{\gamma_1}{k + \gamma_2}\right) = \eta_k$

## □ Example 1: LMS algorithm

$$\nabla f(\underline{w}) = \sum_{l=1}^m \nabla f_l(\underline{w}) = -\sum_{l=1}^m e(l) \underline{p}(l) \Rightarrow \underline{w}(l+1) = \underline{w}(l) + \eta \rho(l) \underline{p}(l)$$

- LMS = incremental SD



# Incremental Gradient Algorithm -2

$$\begin{aligned}\underline{w}(l+1) &= \underline{w}(l) + \eta [d(l) - \underline{p}^T(l)\underline{w}(l)] \underline{p}(l) \\ &= [I - \eta \underline{p}(l) \underline{p}^T(l)] \underline{w}(l) + \eta d(l) \underline{p}(l)\end{aligned}$$

- Convergence for  $0 < \eta < \frac{2}{\lambda_{\max}}$ ;  $\lambda_{\max} = \lambda_{\max} [ \frac{1}{l} \sum \underline{p}(l) \underline{p}^T(l) ] = R$

$$0 < \eta < \frac{2}{\text{tr}(R)} = \frac{2}{\frac{1}{l} \sum p^T(l) p(l)}$$

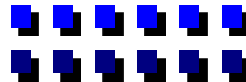
## □ Example 2: Back Propagation Neural Networks

$$w_{ji}^q(l+1) = w_{ji}^q(l) + \underbrace{\eta \delta_j^q(l) y_i^{(q-1)}(l)}_{\text{delta rule } \Delta w_{ji}^q(l)}$$

- Steepest descent with constant step size...very slow.
- Usually employ heavy ball method...adding momentum term to the update.

$$\underline{\Psi}_l = \underline{\Psi}_{l-1} + \eta_k \nabla \underline{g}_i(\underline{x}_k) g_i(\underline{x}_k) + \beta (\underline{\Psi}_{l-1} - \underline{\Psi}_{l-2})$$

$$\text{NN: } w_{ji}^q(l+1) = w_{ji}^q(l) + \underbrace{\eta \delta_j^q(l) y_i^{(q-1)}(l)}_{\text{delta rule}} + \beta (w_{ji}^q(l) - w_{ji}^q(l-1))$$





# Incremental Gradient Algorithm -3

$$\begin{aligned}\Delta w_{ji}(l) &= \beta \Delta w_{ji}^q(l-1) + \eta \delta_j^q(l) y_i^{(q-1)}(l) \\ &= \eta \sum_{k=0}^l \beta^{l-k} \delta_j^q(k) y_i^{(q-1)}(k) = -\eta \sum_{k=0}^l \beta^{l-k} \frac{\partial f}{\partial w_{ji}(k)} \quad 0 < \beta < 1\end{aligned}$$

$\Rightarrow$  same sign for  $\frac{\partial f}{\partial w_{ji}(k)} \Rightarrow \Delta w_{ji} \uparrow$

$\Rightarrow$  alternate sign for  $\frac{\partial f}{\partial w_{ji}(k)} \Rightarrow \Delta w_{ji} \downarrow$

- Randomize data blocks  $\{l^s\}$
- Use non-symmetric sigmoids.

$$\phi(v) = a \tanh(\sigma v) = \frac{2a}{1 + e^{-\sigma v}} - a$$

$$a = 1.716, \quad \sigma = \frac{2}{3}$$

- Conjugate gradient algorithms:

FR

SW

PRP



# Incremental Gradient Algorithm -4

- Incremental Newton = Incremental Gauss-Newton = RLS

□ Newton and Linear Least Squares:

$$\underline{x}_{k+1} = \underline{x}_k - \left( \sum_{i=1}^m \underline{a}_i \underline{a}_i^T \right)^{-1} \left[ A \underline{b} - \sum_{i=1}^m \underline{a}_i \underline{a}_i^T \underline{x}_k \right] = (A A^T)^{-1} A \underline{b}$$

$$\Psi_0 = \underline{x}_k; \quad P_0 = 10^6 I$$

FOR  $i = 1, \dots, m$  DO

$$\underline{k}_i = \frac{P_{i-1} \underline{a}_i}{1 + \underline{a}_i^T P_{i-1} \underline{a}_i}$$

$$\underline{\Psi}_i = \underline{\Psi}_{i-1} + \underline{k}_i (b_i - \underline{a}_i^T \underline{\Psi}_{i-1})$$

$$P_i = P_{i-1} - \frac{P_{i-1} \underline{a}_i \underline{a}_i^T P_{i-1}}{1 + \underline{a}_i^T P_{i-1} \underline{a}_i}$$

ENDDO

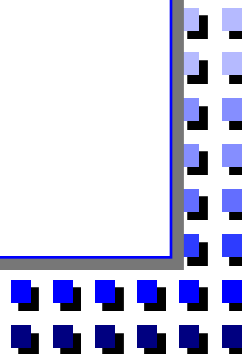
*Sherman – Morrison – Woodbury*

$$\left( \sum_{i=1}^m \underline{a}_i \underline{a}_i^T \right)^{-1} = P_{i-1}$$

$$P_i^{-1} = P_{i-1}^{-1} + \underline{a}_i \underline{a}_i^T$$

$$P_i = P_{i-1} - \frac{P_{i-1} \underline{a}_i \underline{a}_i^T P_{i-1}}{1 + \underline{a}_i^T P_{i-1} \underline{a}_i}$$

$\Rightarrow$  want to put more emphasis on recent data:  $P_{i-1} \rightarrow \frac{P_{i-1}}{\lambda}; 0 < \lambda \leq 1$





# Incremental Gauss-Newton = EKF -1

- Idea can be extended to nonlinear least squares:

$$\underline{x}_{k+1} = \arg \min_{\underline{x}} \frac{1}{2} \sum_{i=1}^m \left\| g_i(\underline{x}_k) + \nabla \underline{g}_i^T(\underline{x}_k)(\underline{x} - \underline{x}_k) \right\|_2^2$$

$$\underline{x}_{k+1} = \underline{x}_k - \left[ \sum_{i=1}^m \nabla \underline{g}_i(\underline{x}_k) \nabla \underline{g}_i^T(\underline{x}_k) \right]^{-1} \sum_{i=1}^m g_i(\underline{x}_k) \nabla \underline{g}_i(\underline{x}_k)$$

- If we want to implement this recursively, one way to do this is as follows:

$$\Psi_0 = \underline{x}_k; \quad P_0 = 10^6 I$$

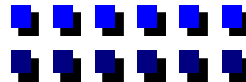
FOR  $i = 1, \dots, m$  DO

$$\underline{k}_i = \frac{P_{i-1} \nabla \underline{g}_i(\underline{x}_k)}{1 + \nabla \underline{g}_i^T(\underline{x}_k) P_{i-1} \nabla \underline{g}_i(\underline{x}_k)}$$

$$\underline{\Psi}_i = \underline{\Psi}_{i-1} + \underline{k}_i g_i(\underline{\Psi}_{i-1}) \approx \underline{\Psi}_{i-1} + \underline{k}_i \left[ g_i(\underline{x}_k) + \nabla \underline{g}_i^T(\underline{x}_k)(\underline{\Psi}_{i-1} - \underline{x}_k) \right]$$

$$P_i = P_{i-1} - \frac{P_{i-1} \nabla \underline{g}_i(\underline{x}_k) \nabla \underline{g}_i^T(\underline{x}_k) P_{i-1}}{1 + \nabla \underline{g}_i^T(\underline{x}_k) P_{i-1} \nabla \underline{g}_i(\underline{x}_k)}$$

END





# Incremental Gauss-Newton = EKF -2

$\Rightarrow$  Forgetting factor  $P_{i-1} \rightarrow P_{i-1}/\lambda$   $0 < \lambda \leq 1$

This is a special case of EKF.

## □ Application to Neural Networks

$$d(i) = h(\underline{w}^T(i) \underline{y}(i)) + e(i) \quad \text{error}$$

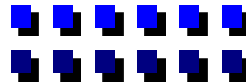
Suppose we have estimate  $\hat{\underline{w}}(i)$ :

$$d(i) \approx h(\hat{\underline{w}}^T(i) \underline{y}(i)) + \nabla h_{\underline{w}}^T(\hat{\underline{w}}) \underbrace{[\underline{w}(i) - \hat{\underline{w}}(i)]}_{\Delta \underline{w}(i)} + e(i)$$

$$\Delta d(i) = d(i) - h(\hat{\underline{w}}^T(i) \underline{y}(i))$$

$$\Delta d(i) = \underline{q}^T(i) \Delta \underline{w}(i) + e(i)$$

$$\underline{q}(i) = \sigma h(\hat{\underline{w}}^T(i) \underline{y}(i)) \left[ 1 - h(\hat{\underline{w}}^T(i) \underline{y}(i)) \right] \underline{y}(i)$$





## Incremental Gauss-Newton = EKF -2

- Application to Neural Networks, continued

$$\Psi_0 = \underline{\hat{w}}(0); \quad P_0 = 10^6 I$$

FOR  $i = 1, \dots, m$  DO

$$\underline{k}_i = \frac{P_{i-1} \underline{q}(i)}{1 + \underline{q}^T(i) P_{i-1} \underline{q}(i)}$$

$$\underline{\hat{w}}_i = \underline{\hat{w}}_{i-1} + \underline{k}_i \left[ d(i) - h(\underline{\hat{w}}_{i-1} \underline{y}(i)) + \underline{q}_i^T (\underline{\hat{w}}_{i-1} - \underline{\hat{w}}_{i-2}) \right]$$

$$P_i = P_{i-1} - \frac{P_{i-1} \underline{q}(i) \underline{q}^T(i) P_{i-1}}{1 + \underline{q}^T(i) P_{i-1} \underline{q}(i)}$$

END

$$\text{Fading memory } P_{i-1} \rightarrow \frac{P_{i-1}}{\lambda} \quad 0 < \lambda \leq 1$$





# Summary

- ❑ Quasi-Newton (QN) Methods
- ❑ Motivation from Quadratic Functions
- ❑ Square-root Implementation
  - DFP & BFGS versions
  - Relation to Filtering with Perfect Observations
- ❑ Properties of QN Methods
  - Quadratic Termination
  - Convergence Properties
  - Scaling
- ❑ Incremental Methods for Neural Networks
  - Back propagation, momentum methods
  - LMS and Recursive least squares
  - Gauss-Newton  $\Rightarrow$  Extended Kalman Filter