

Optimization

Find local minimum of $f(x)$

Brent's Method \rightarrow bracketing

Newton's Method: find $f'(x) = 0$

let x_n be a tentative value near the minimum x^* .

Taylor Series

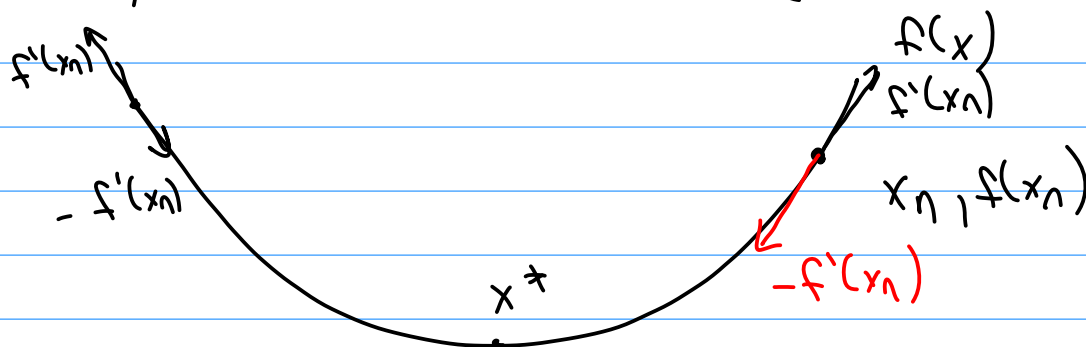
$$f'(x) = f'(x_n) + f''(x_n)(x - x_n) + \frac{1}{2} f'''(x_n)(x - x_n)^2 + \text{H.O.T.}$$

$$\hat{f}'(x) = f'(x_n) + f''(x_n)(x - x_n) = 0$$

$$\Rightarrow x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

- 1) might diverge if $f''(x_n) \sim 0$
- 2) might oscillate

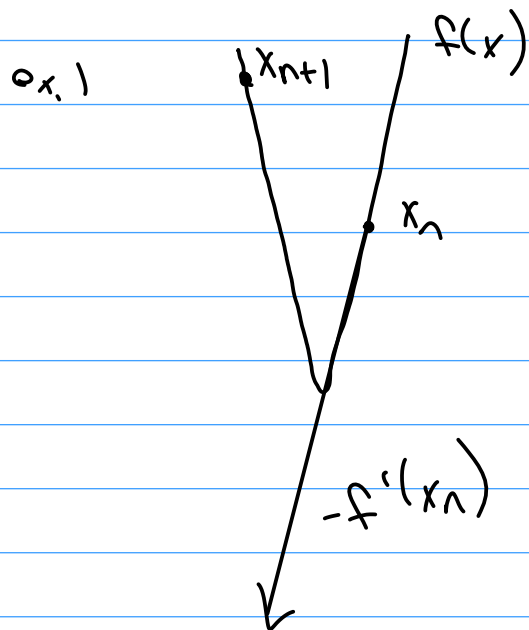
Steepest Gradient Descent



$$x_{n+1} = x_n - \alpha_n f'(x_n)$$

\uparrow will need damping

You need damping as you do not know how $f'(x)$ behaves



Multivariable Minimization

$$\text{Minimize } f(\underline{x}) \quad \text{w/ } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

\uparrow
 A scalar

ex1) minimize $f(x_1, x_2) = \sin(x_1) \cos(x_2)$

Minimum occurs where

$$\nabla f(\underline{x}^*) = \begin{bmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Taylor Series for \underline{x}_n near minimum, \underline{x}^*

$$\nabla f(\underline{x}) \approx \nabla f(\underline{x}_n) + \underline{H}(\underline{x}_n)(\underline{x} - \underline{x}_n) + \text{H.O.T.}$$

$\underline{H}(\underline{x}) =$ Hessian Matrix of $f(\underline{x})$

$$= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & & & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

$$\underline{H}(\underline{x}) = \underline{H}^T(\underline{x})$$

To first-order, $\hat{\nabla} f(\underline{x}) = \nabla f(\underline{x}_n) + \underline{H}(\underline{x}_n)(\underline{x} - \underline{x}_n)$

$$\hat{\nabla} f(\underline{x}_{n+1}) = \underline{0}$$

$$\Rightarrow \underline{x}_{n+1} = \underline{x}_n - \underline{H}_n^{-1} \nabla f(\underline{x}_n)$$

Write as:

1) Solve $\underline{H}_n \underline{d}_n = -\nabla f_n$

2) Update: $\underline{x}_{n+1} = \underline{x}_n + \alpha_n \underline{d}_n$

Notes: - \underline{H}_n changes every iteration

- Need to solve a linear system each iteration

- \underline{H}_n might be expensive or unknown

Quasi-Newton Methods

Idea: Find a $B_n \in H_n$ where
 B_n is cheap to compute &
solving $B_n d_n = -\nabla f_n$ is also cheap

Need: A sequence B_0, B_1, B_2, \dots

1) How to choose B_0 ?

2) How does B_n depend on $B_{n-1}, B_{n-2},$
solution, etc.

General Quasi-Newton Algorithm

let x_0 & B_0 be known

$n=0$

until converged

Solve $B_n d_n = -\nabla f_n$

$x_{n+1} = x_n + \alpha_n d_n$ via line search

$y_n = \nabla f(x_{n+1}) - \nabla f(x_n) = \nabla (f(x_{n+1}) - f(x_n))$

$B_{n+1} = \text{function of } B_n, d_n, y_n$

Simple choice for B_0

let $\hat{d}_0 = -\nabla f(x_0)$ $y_0 = \nabla (f(x_0 + \hat{d}_0) - f(x_0))$

$$B_0 = \frac{y_0^T y_0}{y_0^T \hat{d}_0} I \quad B_0^{-1} = \frac{y_0^T \hat{d}_0}{y_0^T y_0} I$$

Thus, in iteration $d_0 = -\frac{y_0^T \hat{d}_0}{y_0^T y_0} \nabla f_0$

Many choices for \underline{B}_n

- Davidon - Fletcher - Powell (DFP) Method

$$\underline{B}_{n+1} = \left(\underline{I} - \frac{\underline{y}_n \underline{d}_n^T}{\underline{y}_n^T \underline{d}_n} \right) \underline{B}_n \left(\underline{I} - \frac{\underline{d}_n \underline{y}_n^T}{\underline{y}_n^T \underline{d}_n} \right) + \frac{\underline{y}_n \underline{y}_n^T}{\underline{y}_n^T \underline{y}_n}$$

$$\underline{B}_{n+1}^{-1} = \underline{B}_n^{-1} + \frac{\underline{d}_n \underline{d}_n^T}{\underline{d}_n^T \underline{y}_n} - \underline{B}_n^{-1} \left(\frac{\underline{y}_n \underline{y}_n^T}{\underline{y}_n^T \underline{B}_n^{-1} \underline{y}_n} \right) \underline{B}_n^{-1}$$

(Given \underline{B}_0^{-1} , \underline{B}_1^{-1} is known,
(Given \underline{B}_1^{-1} , \underline{B}_2^{-1} is known, etc.

- Broyden - Fletcher - Goldfarb - Shanno (BF(GS))

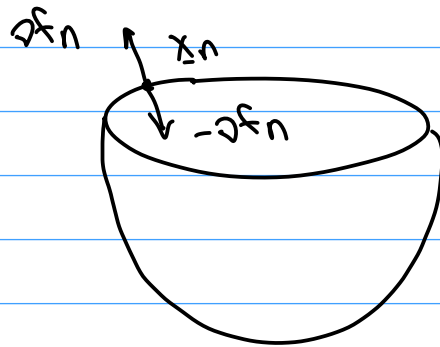
Based on outer products.

$$\underline{B}_{n+1} = \underline{B}_n + \frac{\underline{y}_n \underline{y}_n^T}{\underline{y}_n^T \underline{d}_n} - \underline{B}_n \left(\frac{\underline{d}_n \underline{d}_n^T}{\underline{d}_n^T \underline{B}_n \underline{d}_n} \right) \underline{B}_n^T$$

$$\underline{B}_{n+1}^{-1} = \left(\underline{I} - \frac{\underline{d}_n \underline{y}_n^T}{\underline{y}_n^T \underline{d}_n} \right) \underline{B}_n^{-1} \left(\underline{I} - \frac{\underline{y}_n \underline{d}_n^T}{\underline{y}_n^T \underline{d}_n} \right) + \frac{\underline{d}_n \underline{y}_n^T}{\underline{y}_n^T \underline{d}_n}$$

Multidimensional Gradient Descent

A 2D convex function $f(x_1, x_2)$



$$\Rightarrow \underline{x}_{n+1} = \underline{x}_n - \alpha_n \nabla f(\underline{x}_n)$$

Line Search / Backtracking

Many methods written as $\underline{x}_{n+1} = \underline{x}_n + \alpha_n \underline{d}_n$

Newton: $\underline{d}_n = -\underline{J}_n^{-1} \underline{f}_n$

(Gradient Descent: $\underline{d}_n = -\nabla f(\underline{x}_n)$)

Minimization: $\underline{d}_n = -\underline{H}_n^{-1} \nabla f(\underline{x}_n)$

In all cases you want

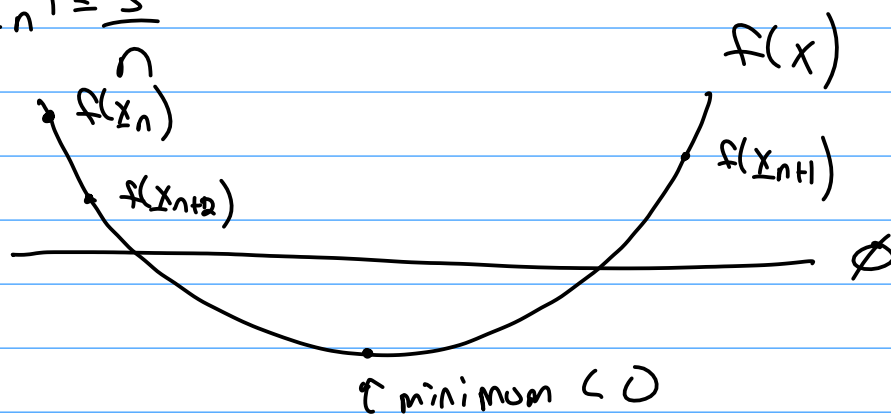
$f(\underline{x}_n + \alpha_n \underline{d}_n) < f(\underline{x}_n)$ minimization

$\|\underline{f}(\underline{x}_n + \alpha \underline{d}_n)\| < \|\underline{f}(\underline{x}_n)\|$ root finding

Focus on $f(\underline{x})$ (or $f(\underline{x}) = \|\underline{g}(\underline{x})\|$)

It is not sufficient to simply
require $f(\underline{x}_n + \alpha_n \underline{d}_n) < f(\underline{x}_n) \leftarrow$

ex) Some minimization problem has
 $f(\underline{x}_n) = \underline{5}$



minimum is never achieved.

Wolfe Conditions

Require that

$$1) \quad f(\underline{x}_n + \alpha_n \underline{d}_n) \leq f(\underline{x}_n) + C_1 \alpha_n (\nabla f(\underline{x}_n))^T \underline{d}_n \\ \text{for some } C_1 \in (0, 1)$$

This requires **sufficient progress** towards the solution,

In practice C_1 is small, say $C_1 \approx 10^{-4}$

2) To avoid small α_n require that

$$(\nabla f(\underline{x}_n + \alpha_n \underline{d}_n))^T \underline{d}_n \geq C_2 (\nabla f(\underline{x}_n))^T \underline{d}_n \\ \text{for some } C_2 \in (C_1, 1)$$

Thus $0 < C_1 < C_2 < 1$

It can be shown that if $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable, \underline{d}_n is a descent direction at \underline{x}_n , and if $0 < C_1 < C_2 < 1$ there exists a range of α_n that satisfy the Wolfe Conditions

One method for α_n :

$$\alpha_n = \frac{(\overbrace{\underline{x}_n - \underline{x}_{n-1}}^{\underline{d}_{n-1}})^T (\overbrace{\nabla f(\underline{x}_n) - \nabla f(\underline{x}_{n-1})}^{\underline{y}_{n-1}})}{\|\nabla f(\underline{x}_n) - \nabla f(\underline{x}_{n-1})\|_2^2}$$

Another method

Define $\phi(\alpha) = f(\underbrace{x_n}_{\text{known}} + \alpha \underline{d}_n)$ $\phi(0) = f(x_n)$

$$\phi'(\alpha) = \underline{d}_n \cdot \nabla f(x_n + \alpha \underline{d}_n) \quad \phi'(0) = \underline{d}_n \cdot \nabla f(x_n)$$

Wolfe Condition #1: $\phi(\alpha) \leq \phi(0) + C_1 \alpha \phi'(0)$

(General Idea: choose α_0 , say $\alpha_0 = 1$ (full step))

If $\phi(\alpha_0) \leq \phi(0) + C_1 \alpha_0 \phi'(0)$ use α_0

If not satisfied actual $\alpha \in (0, \alpha_0)$

Currently (given: $\phi(0), \phi'(0), \phi(\alpha_0)$)

Form an interpolant

$$\hat{\phi}(\alpha) = \left(\frac{\phi(\alpha_0) - \phi(0) - \alpha_0 \phi'(0)}{\alpha_0^2} \right) \alpha^2 + \phi'(0) \alpha + \phi(0)$$

Find minimum of $\hat{\phi}(\alpha)$

$$\alpha_1 = \frac{-\phi'(0) \alpha_0^2}{2[\phi(\alpha_0) - \phi(0) - \phi'(0) \alpha_0]}$$

If $\phi(\alpha_1) \leq \phi(0) + C_1 \alpha_1 \phi'(0)$ use α_1

If not form a cubic polynomial w/
 $\phi(u), \phi'(u), \phi(\alpha_0), \phi(\alpha_1)$

$$\tilde{\phi}(\alpha) = a\alpha^3 + b\alpha^2 + c\alpha + d$$

$$\text{Solve } \tilde{\phi}(\alpha_2) = 0$$

$$\alpha_2 = \frac{-b + \sqrt{b^2 - 3a\phi'(u)}}{3a}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \frac{1}{\alpha_0^2 \alpha_1^2 (\alpha_1 - \alpha_0)} \begin{bmatrix} \alpha_0^2 & -\alpha_1^2 \\ -\alpha_0^3 & \alpha_1^3 \end{bmatrix} \begin{bmatrix} \phi(\alpha_1) - \phi(u) - \phi'(u)\alpha_1 \\ \phi(\alpha_0) - \phi(u) - \phi'(u)\alpha_0 \end{bmatrix}$$

check if $\phi(\alpha_n) \leq \phi(u) + C_1 \alpha_n \phi'(u)$

If not true repeat $\phi(u), \phi'(u), \phi(\alpha_{k-1}), \phi(\alpha_{k-2})$

Sequence of

$$\alpha_0 > \alpha_1 > \alpha_2 > \dots > \alpha_{k-1} > \alpha_k > \dots > 0$$

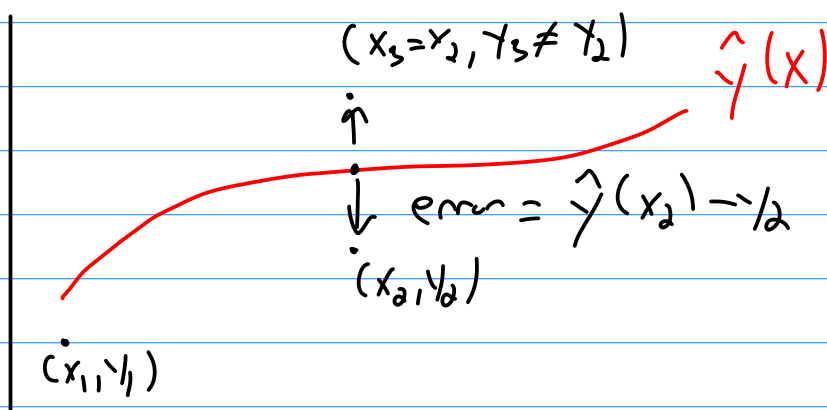
If any $\alpha_k - \alpha_{k-1} < \varepsilon$ or $\alpha_k \ll \alpha_{k-1}$

$$\text{or } \alpha_k = \frac{1}{2} \alpha_{k-1}$$

See "Numerical Optimization" by
Nocedal & Wright

Regression : Curve fitting

Regression \rightarrow fit an equation to a set of Data but do not require that it pass through the Data.



(Goal: minimize the error

Define the error @ point (x_i, y_i) as
$$r_i = y_i - \hat{y}(x_i, \underline{a})$$

\underline{a} = list of coefficients in the model equation

ex.) $\hat{y}(x, \underline{a}) = a_0 + a_1 x + a_2 x^2 + a_3 x^3$
example of linear regression

ex.) $\hat{y}(x, \underline{a}) = a_0 \sin(x) + a_1 \sin(2x) + a_2 \cos(x) + a_3 \cos(2x)$

Linear regression
ex.) $\hat{y}(x, \underline{a}) = \frac{a_1 x}{a_2 + x}$ non linear regression

$$\text{ex. 1 } \hat{y}(x, \underline{a}) = a_0 \sin(a_1 x) + a_2 \cos(a_2 x)$$

non-linear regression

For both, need to minimize $\|\underline{r}\|_2$

$$\underline{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix} = \begin{bmatrix} y_1 - \hat{y}(x_1, \underline{a}) \\ y_2 - \hat{y}(x_2, \underline{a}) \\ \vdots \\ y_n - \hat{y}(x_n, \underline{a}) \end{bmatrix} \quad \begin{array}{l} \text{for } n\text{-pts} \\ \text{Points} \end{array}$$

For Linear regression:

$$\hat{y}(x, \underline{a}) = a_1 f_1(x) + a_2 f_2(x) + \dots + a_p f_p(x)$$

$$\hat{y}(x_1, \underline{a}) = a_1 f_1(x_1) + a_2 f_2(x_1) + \dots + a_p f_p(x_1) = y_1$$

$$\hat{y}(x_2, \underline{a}) = a_1 f_1(x_2) + \dots + a_p f_p(x_2) = y_2$$

$$\hat{y}(x_n, \underline{a}) = a_1 f_1(x_n) + \dots + a_p f_p(x_n) = y_n$$

\Downarrow

$$\begin{bmatrix} f_1(x_1) & f_2(x_1) & \dots & f_p(x_1) \\ f_1(x_2) & f_2(x_2) & \dots & f_p(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ f_1(x_n) & f_2(x_n) & \dots & f_p(x_n) \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_p \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$\begin{matrix} \text{Data} \\ \text{points} \end{matrix}$

$\begin{matrix} \text{Functions} \\ \text{Data} \end{matrix}$

$\begin{matrix} \text{Parameters} \\ \text{Data} \end{matrix}$

In general $n > p \Rightarrow \text{minimize } \|\underline{r}\|_2$

The **Normal Equations** minimize $\|\underline{r}\|_2$

(Given $\underline{F}\underline{q} = \underline{y}$, the solution to

$$\underline{F}^T \underline{F} \hat{\underline{q}} = \underline{F}^T \underline{y} \quad \text{minimizes}$$

$$\|\underline{F}\hat{\underline{q}} - \underline{y}\|_2$$

Also called the **least squares solution**

$$\hat{\underline{q}} = (\underline{F}^T \underline{F})^{-1} \underline{F}^T \underline{y} \quad \text{is unique}$$

Never actually form the Normal Equation, why?

Condition Number

$$K(\underline{A}^T) = K(\underline{A})$$

$$K(\underline{A}\underline{B}) = K(\underline{A})K(\underline{B})$$

$$K(\underline{A}^T \underline{A}) = K(\underline{A}^T) K(\underline{A}) = (K(\underline{A}))^2$$

but, other methods can find the least squares solution (QR, SVD, etc.)

In Matlab: Backslash $\underline{q} = \underline{F} \setminus \underline{y}$

Nonlinear: No simple linear system

Define a objective function $S(\underline{q}) = \|\underline{r}\|_2^2$

$$S(\underline{a}) = \sum_{i=1}^n r_i^2 \quad \text{Use any minimization method.}$$

look at Gauss-Newton

(The information below has been expanded on past that presented in lecture)

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

$$\underline{g}_k = \nabla S(\underline{a}_k) \rightarrow g_j = \frac{\partial S}{\partial a_j}$$

$$\underline{H}_k = \text{Hessian of } S \text{ @ } \underline{a}_k$$

$$\Rightarrow H_{j,k} = \frac{\partial^2 S}{\partial a_j \partial a_k} = \frac{d}{da_k} \left(\frac{dS}{da_j} \right) = \frac{dg_j}{da_k}$$

$$S(\underline{a}) = \sum_{i=1}^n (y_i - \hat{y}(x_i, \underline{a}))^2 = \sum_{i=1}^n r_i^2$$

$$g_j = \frac{d}{da_j} \left(\sum_{i=1}^n r_i^2 \right) = 2 \sum_{i=1}^n r_i \frac{\partial r_i}{\partial a_j} = 2 \sum_{i=1}^n r_i \underset{\substack{\uparrow \\ \text{Jacobian of } r(\underline{a})}}{J_{ij}}$$

$$\Rightarrow \underline{g} = 2 \underline{r}^T \underline{J}$$

$$H_{j,k} = 2 \sum_{i=1}^n \left(\underbrace{\frac{\partial r_i}{\partial a_j} \frac{\partial r_i}{\partial a_k}}_{1^{\text{st}} \text{-order}} + \underbrace{r_i \frac{\partial^2 r_i}{\partial a_j \partial a_k}}_{2^{\text{nd}} \text{-order}} \right)$$

neglect 2nd-order part

$$\text{Thus } H_{jk} \approx 2 \sum_{i=1}^n \frac{\partial r_i}{\partial a_j} \frac{\partial r_i}{\partial a_k} = 2 \sum_{i=1}^n J_{ij} J_{ik}$$

$$\Rightarrow \underline{H} = 2 \underline{J}^T \underline{J}$$

$$\Rightarrow \underline{a}_{k+1} = \underline{a}_k - \underline{H}_k^{-1} \underline{g}_k \approx \underline{a}_k - \frac{1}{2} (\underline{J}_k^T \underline{J}_k)^{-1} (2 \underline{J}_k^T \underline{r}_k)$$

$$\underline{a}_{k+1} = \underline{a}_k - (\underline{J}_k^T \underline{J}_k)^{-1} (\underline{J}_k^T \underline{r}_k)$$

$$\text{Define } \underline{d}_k = \underline{a}_{k+1} - \underline{a}_k$$

$$\Rightarrow \underline{d}_k = -(\underline{J}_k^T \underline{J}_k)^{-1} (\underline{J}_k^T \underline{r}_k)$$

$$\underline{J}_k^T \underline{J}_k \underline{d}_k = -\underline{J}_k^T \underline{r}_k \leftarrow \text{normal equations for } \underline{J} \underline{d}_k = -\underline{r}_k!$$

Thus, solve $\underline{J}_k \underline{d}_k = -\underline{r}_k$ in least-squares sense, then $\underline{a}_{k+1} = \underline{a}_k + \underline{d}_k$

Advantage: No Hessian!

Disadvantage: Only converges if

$$\left| r_i \frac{\partial^2 r_i}{\partial a_j \partial a_k} \right| \ll \left| \frac{\partial r_i}{\partial a_j} \frac{\partial r_i}{\partial a_k} \right|$$

\uparrow \uparrow

Part that is
ignored

Part that is
kept

Typically valid if $\|h\|$ is initially
small or if $\hat{y}(x, \underline{a})$ is
mildly non-linear (no large
oscillations).