

# Optimization

Find local minimum of  $f(x)$

Brent's Method → bracketing

---

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

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

Taylor Series

$$f(x) = f'(x_n) + f''(x_n)(x-x_n)$$

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

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

$\Rightarrow$

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

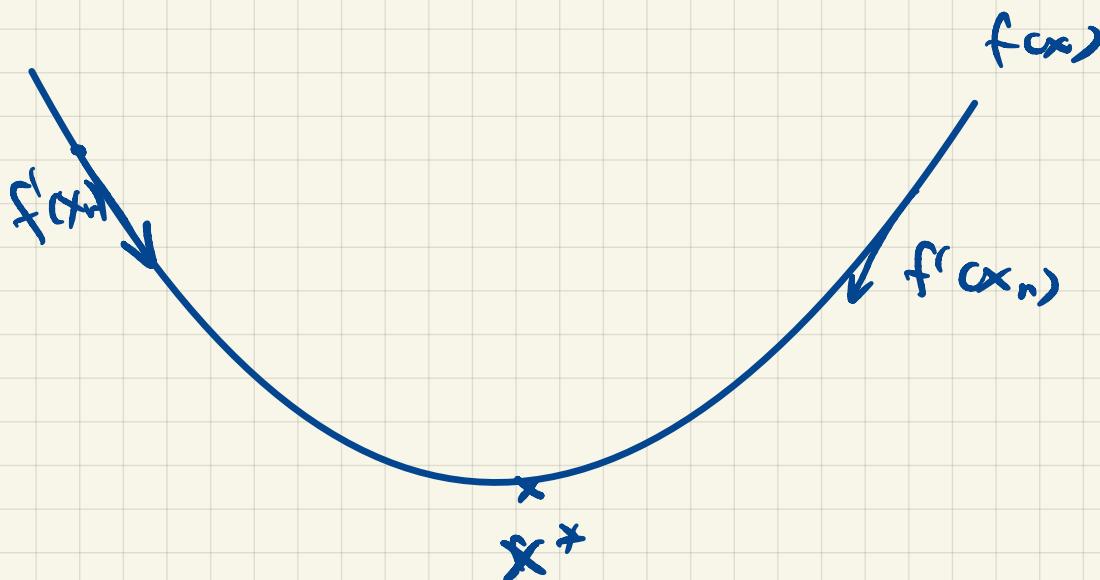
① might diverge  $\Rightarrow f''(x_n) \sim 0$

② might oscillate

Steepest

Gradient

Descent

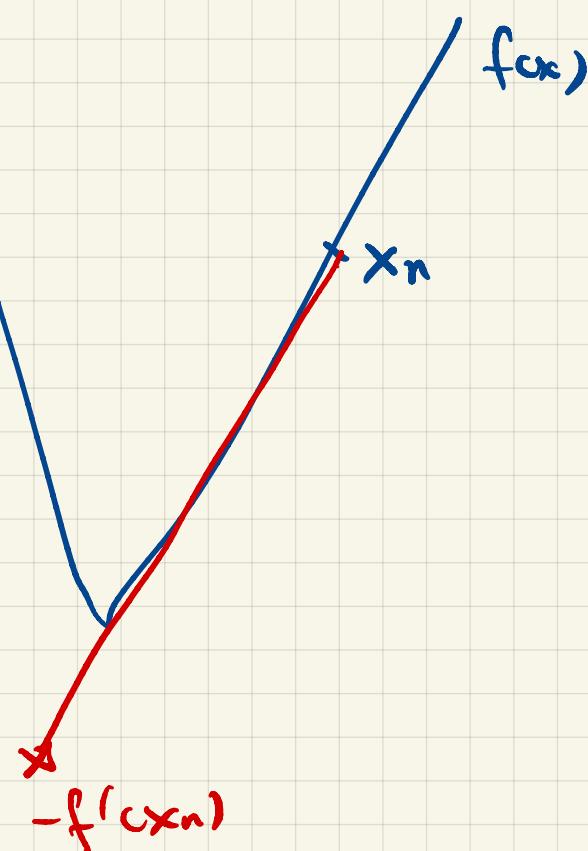


$$x_{n+1} = x_n - \alpha_n f'(x_n)$$

↑  
need damping

need damping

because the  
behavior of  $f(x)$   
is unclear



## Multi variable

## Minimization

Minimize  $f(\underline{x})$   
a scalar

$$\text{w/ } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

Minimum occurs where

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

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

$$\nabla f(\underline{x}) = \nabla f(\underline{x}_n) + H(\underline{x}_n)(\underline{x} - \underline{x}_n) + \underline{\text{higher order terms}}$$

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

$$H(\underline{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & & & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

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

To first order

$$\nabla \hat{f}(\underline{x}_n) = \nabla f(x_n) + H(x_n)(\underline{x} - x_n)$$

$$\nabla \hat{f}(x_{n+1}) = 0$$

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

$$\textcircled{1} \quad \text{Solve} \quad \underline{H}_n \begin{matrix} \underline{s}_n \\ \downarrow \\ (\underline{x}_{n+1} - \underline{x}_n) \end{matrix} = -\nabla f_n$$

$$\textcircled{2} \quad \text{Update} \quad \underline{x}_{n+1} = \underline{x}_n + \alpha_n \underline{s}_n$$

Note:

- ①  $\underline{H}_n$  changes every iteration
- ② Need to solve a linear system each iteration
- ③  $\underline{H}_n$  might be expensive or unknown

## Quasi-Newton Method

Idea: Find a  $\underline{B_n} \approx \underline{H_n}$  where

$\underline{B_n}$  is cheap to compute and

solving  $\underline{B_n} \underline{s_n} = -\nabla \underline{f_n}$  is also cheap

Need: A sequence  $\underline{B_0}, \underline{B_1}, \underline{B_2}, \dots$

1) How to choose  $\underline{B_0}$ ?

2) How does  $\underline{B_n}$  update based  
on  $\underline{B_{n-1}}, \underline{B_{n-2}}, \dots$

## General Quasi-Newton Algorithm

let  $x_0$  &  $\underline{B_0}$  be known

$n=0$

until converge

Solve  $\underline{B_n} \underline{s_n} = -\nabla \underline{f_n}$

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

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

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

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

---

Choice for  $B_0$

Let  $\tilde{s}_0 = -\nabla f(x_0)$

$$y_0 = \nabla(f(x + \delta_0) - f(x))$$

$$B_0 = \frac{\tilde{s}_0^T \tilde{s}_0}{\tilde{s}_0^T y_0} =$$

$$B_0^{-1} = \frac{\tilde{s}_0^T \delta_0}{\tilde{s}_0^T y_0} =$$

Then iteration,

$$\delta_0 = -\frac{\tilde{s}_0^T \delta_0}{\tilde{s}_0^T y_0} \nabla f_0$$

Many Choices for  $\underline{B}_n$

- DFP method

Davidson - Fletcher - Powell

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

$$+ \frac{\underline{y}_n \underline{y}_n^T}{\underline{y}_n^T \underline{y}_n}$$

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

- BFGS method

Broyden - Fletcher - Goldfarb

- Shanno

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

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

$$+ \frac{\underline{s}_n \underline{y}_n^\top}{\underline{y}_n^\top \underline{s}_n}$$


---

Line Search / Backtracking

$$\underline{x}_{n+1} = \underline{x}_n + \alpha_n \underline{s}_n$$

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

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

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

In all cases,

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

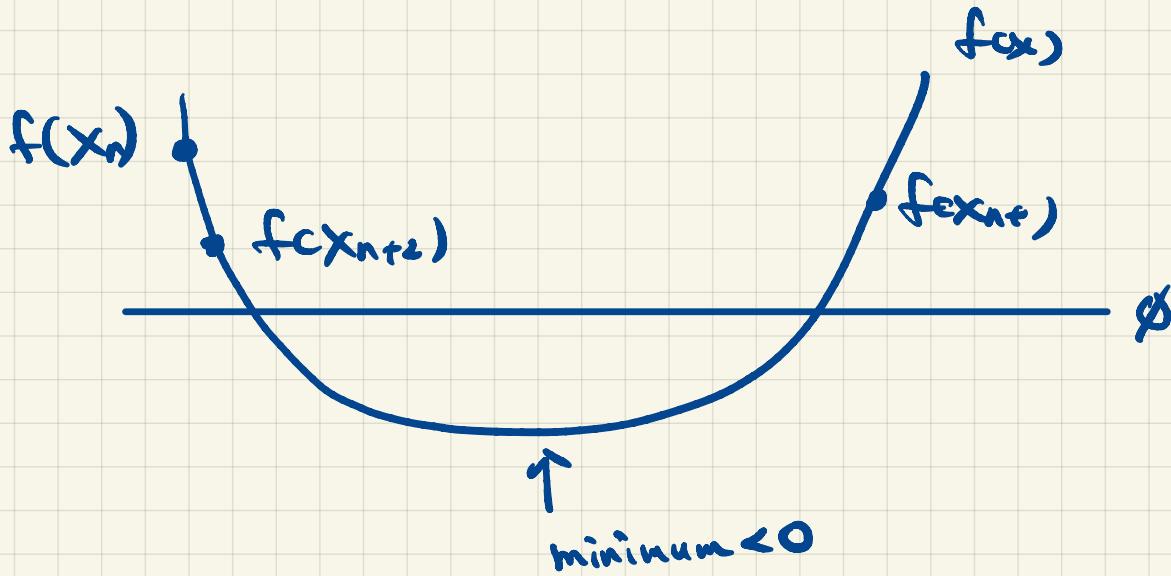
minimization

$$\|f(\underline{x}_n + \alpha_n \underline{s}_n)\| < \|f(\underline{x}_n)\|$$

Focus on  $f(x)$ )

It is not sufficient to

simply require



Wolfe Conditions.

$$\textcircled{1} \quad f(x_n + \alpha_n \delta_n) \leq f(x_n) + c_1 \alpha_n (\nabla f_n)^T \delta_n$$

for some  $c_1 \in [0, 1]$

This requires sufficient progress towards the sol.

$$c_1 \sim 10^{-4} \quad \text{small}$$

② To avoid small  $\alpha_n$ , require

$$(\nabla f(\underline{x}_n + \alpha_n \underline{g}_n))^T \underline{g}_n \geq c_2 (\nabla f_n)^T \underline{g}_n$$

for  $c_2 \in (c_1, 1)$

Thus  $0 < c_1 < c_2 < 1$

It can be shown that if

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously

differentiable,  $\underline{g}_n$  is a

descent direction at  $\underline{x}_n$  and if

range of  $\alpha_n$  that satisfy the

Wolfe condition.

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



# 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  $x^*$ .

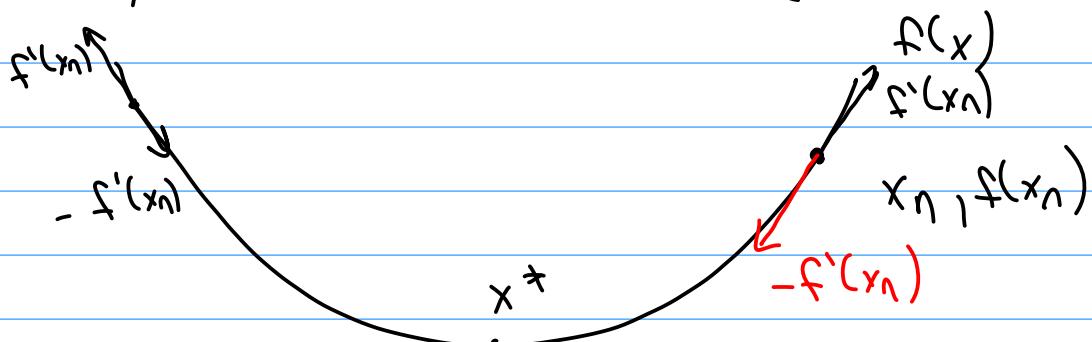
Taylor Series

$$f'(x) = f'(x_n) + f''(x_n)(x - x_n) + \frac{1}{2} f'''(x_n)(x - x_n)^2 + \text{H.O.T.}$$
$$\hat{f}'(x) = f'(x_n) + f''(x_n)(x - x_n) = 0$$

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

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

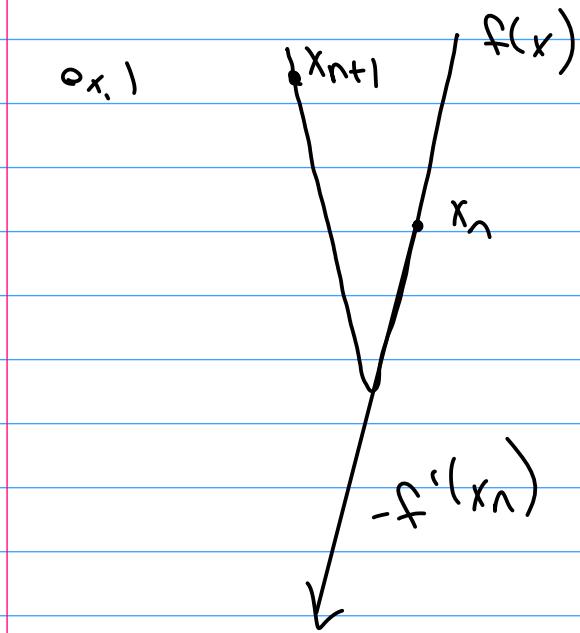
Steepest Gradient Descent



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

↑ will need damping

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



### Multivariable Minimization

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

A scalar

$$\text{ex: } \text{minimize } f(x_1, x_2) = \sin(x_1) \cos(x_2)$$

minimum occurs where

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

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

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

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

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

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

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:

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

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

Note:-  $\underline{H}_n$  changes every iteration

- Need to solve a linear system each iteration
- $\underline{H}_n$  might be expensive or unknown

## Quasi-Newton Methods

Ideas: 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

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

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

---

Simple choice for  $B_0$

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

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

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

Many choices for  $B_n$

- Davidon - Fletcher - Powell (DFP) Method

$$B_{n+1} = \left( I - \frac{y_n d_n^T}{y_n^T d_n} \right) B_n \left( I - \frac{d_n y_n^T}{y_n^T d_n} \right) + \frac{y_n y_n^T}{y_n^T y_n}$$

$$B_n^{-1} = B_0^{-1} + \frac{d_n d_n^T}{d_n^T y_n} - B_0^{-1} \left( \frac{y_n y_n^T}{y_n^T B_0^{-1} y_n} \right) B_0^{-1}$$

(Given  $B_0^{-1}$ ,  $B_0^{-1}$  is known,

(Given  $B_1^{-1}$ ,  $B_1^{-1}$  is known, etc.

- Broyden - Fletcher - Goldfarb - Shanno (BFGS)

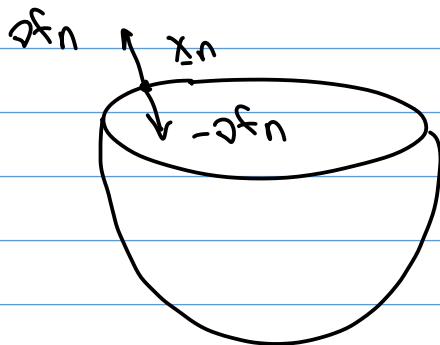
Based on outer products:

$$B_{n+1} = B_n + \frac{y_n y_n^T}{y_n^T d_n} - B_n \left( \frac{d_n d_n^T}{d_n^T B_n d_n} \right) B_n^T$$

$$B_n^{-1} = \left( I - \frac{d_n y_n^T}{y_n^T d_n} \right) B_0^{-1} \left( I - \frac{y_n d_n^T}{y_n^T d_n} \right) + \frac{d_n y_n^T}{y_n^T d_n}$$

## 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  $\underline{x}_{n+1} = \underline{x}_n + \alpha_n \underline{d}_n$

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

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

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

In all cases you want

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

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

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

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

$$f(\underline{x}_n) = \sum$$



minimum is never achieved.

## Wolfe Conditions

Require that

$$1) \quad f(\underline{x}_n + \alpha_n \underline{d}_n) \leq f(\underline{x}_n) + C_1 \alpha_n (\nabla f_n)^T \underline{d}_n$$

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  $\alpha_n$  require that

$$(\nabla f(\underline{x}_n + \alpha_n \underline{d}_n))^T \underline{d}_n \geq C_2 (\nabla f_n)^T \underline{d}_n$$

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  $\alpha_n$  that satisfy the Wolfe conditions

One method for  $\alpha_n$ :

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

Another method

Define  $\phi(\alpha) = f(x_n + \alpha d_n)$   $\phi(0) = f(x_n)$

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

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

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

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

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  $\alpha_1$

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

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

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

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

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

$$\text{check if } \phi(\alpha_2) \leq \phi(\zeta) + C_1 \alpha_2 \phi'(\zeta)$$

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

Sequence of

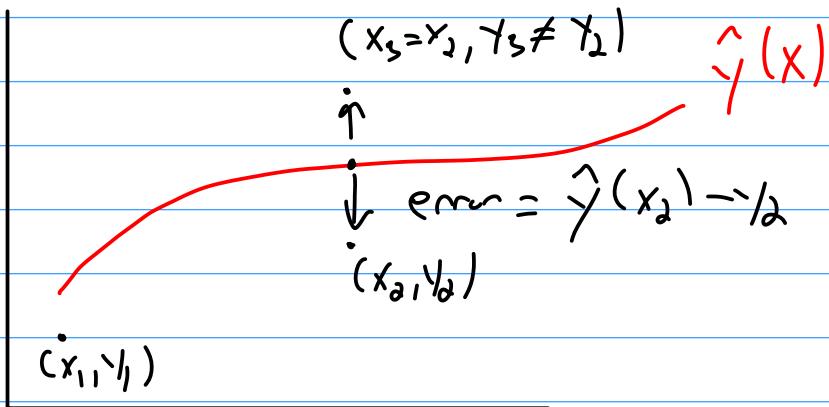
$$\alpha_0 > \alpha_1 > \alpha_2 > \dots > \alpha_{k-1} > \alpha_k > \dots > 0$$

If any  $\alpha_k - \alpha_{k-1} < \varepsilon$  or  $\alpha_k \ll \alpha_{k-1}$   
 $\Rightarrow \alpha_k = \frac{1}{2} \alpha_{k-1}$

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

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

non-linear regression

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

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

For Linear regression:

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

$$\begin{aligned} \hat{y}(x_1, \underline{a}) &= a_1 f_1(x_1) + a_2 f_2(x_1) + \dots + a_p f_p(x_1) = y_1 \\ \hat{y}(x_2, \underline{a}) &= a_1 f_1(x_2) + \dots + a_p f_p(x_2) = y_2 \\ &\vdots \end{aligned}$$

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

$\Downarrow$

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

$E$   
 $n \times p$   
# of data points  
# of functions

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

The Normal Equations minimize  $\|\underline{z}\|_2$

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

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

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

Also called the least squares solution

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

Never actually form the Normal Equations,  
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  $S(\underline{q}) = \|\underline{z}\|_2^2$

$$S(\underline{a}) = \sum_{i=1}^n r_i^2$$

Use any minimization method.

look at Gauss-Newton

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

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

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

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

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

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

$$g_j = \frac{d}{da_j} \left( \sum_{i=1}^n r_i^2 \right) = 2 \sum_{i=1}^n r_i \frac{\partial r_i}{\partial a_j} = 2 \sum_{i=1}^n r_i \underbrace{\mathbf{J}_{ij}}_{\text{Jacobian of } r(\underline{a})}$$

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

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

neglect 2<sup>nd</sup>-order part

$$\text{Thus } H_{jk} \approx 2 \sum_{i=1}^n \frac{\partial r_i}{\partial q_j} \frac{\partial r_i}{\partial q_k} = 2 \sum_{i=1}^n J_{ci} J_{ck}$$

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

$$\Rightarrow \underline{q}_{k+1} = \underline{q}_k - H_k^{-1} \underline{s}_k \approx \underline{q}_k - \frac{1}{2} (\underline{J}_k^T \underline{J}_k)^{-1} (\underline{J}_k^T \underline{r}_k)$$

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

$$\text{Def: } \underline{\delta}_k = \underline{a}_{k+1} - \underline{a}_k$$

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

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

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

Advantage: No Hessian!

Disadvantage: Only convergence if

$$\left| r_i \frac{\partial^2 r_i}{\partial q_i \partial q_k} \right| \ll \left| \frac{\partial r_i}{\partial q_j} \frac{\partial r_i}{\partial q_k} \right|$$

Part that is  
ignored

Part that is  
kept

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