#### Lecture 15

#### Nonlinear Problems

Newton's Method

#### Syllabus

Lecture 01 Describing Inverse Problems Probability and Measurement Error, Part 1 Lecture 02 Probability and Measurement Error, Part 2 Lecture 03 Lecture 04 The L<sub>2</sub> Norm and Simple Least Squares A Priori Information and Weighted Least Squared Lecture 05 **Resolution and Generalized Inverses** Lecture 06 Lecture 07 Backus-Gilbert Inverse and the Trade Off of Resolution and Variance Lecture 08 The Principle of Maximum Likelihood Lecture 09 **Inexact Theories** Lecture 10 Nonuniqueness and Localized Averages Vector Spaces and Singular Value Decomposition Lecture 11 Lecture 12 Equality and Inequality Constraints Lecture 13  $L_1$ ,  $L_{\infty}$  Norm Problems and Linear Programming Lecture 14 Nonlinear Problems: Grid and Monte Carlo Searches **Nonlinear Problems: Newton's Method** Lecture 15 Lecture 16 Nonlinear Problems: Simulated Annealing and Bootstrap Confidence Intervals Lecture 17 **Factor Analysis** Varimax Factors, Empircal Orthogonal Functions Lecture 18 Lecture 19 Backus-Gilbert Theory for Continuous Problems; Radon's Problem Lecture 20 Linear Operators and Their Adjoints Lecture 21 Fréchet Derivatives Lecture 22 Exemplary Inverse Problems, incl. Filter Design Lecture 23 Exemplary Inverse Problems, incl. Earthquake Location Lecture 24 Exemplary Inverse Problems, incl. Vibrational Problems

#### Purpose of the Lecture

#### Introduce Newton's Method

#### Generalize it to an Implicit Theory

Introduce the Gradient Method

#### Part 1

#### Newton's Method

### grid search Monte Carlo Method are completely undirected

### alternative take directions from the *local properties* of the error function *E*(**m**)

#### Newton's Method

#### start with a guess $\mathbf{m}^{(p)}$

near  $\mathbf{m}^{(p)}$ , approximate  $E(\mathbf{m})$  as a parabola and find its minimum

set new guess to this value and iterate



#### Taylor Series Approximation for *E*(**m**)

#### expand *E* around a point $\mathbf{m}^{(p)}$

$$E(\mathbf{m}) \approx E(\mathbf{m}^{(p)}) + \sum_{i=1}^{M} b_i \left( m_i - m_i^{(p)} \right) + \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} B_{ij} \left( m_i - m_i^{(p)} \right) \left( m_j - m_j^{(p)} \right)$$

with 
$$b_i = \frac{\partial E}{\partial m_i} \Big|_{\mathbf{m}^{(p)}}$$
 and  $B_{ij} = \frac{\partial^2 E}{\partial m_i \partial m_j} \Big|_{\mathbf{m}^{(p)}}$ 

#### differentiate and set result to zero to find minimum

$$E(\mathbf{m}) \approx E(\mathbf{m}^{(p)}) + \sum_{i=1}^{M} b_i \left( m_i - m_i^{(p)} \right) + \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} B_{ij} \left( m_i - m_i^{(p)} \right) \left( m_j - m_j^{(p)} \right)$$

$$\frac{\partial E(\mathbf{m})}{\partial m_q} = 0 = b_q + \sum_{j=1}^M B_{qj} \left( m_j - m_j^{(p)} \right)$$

$$\mathbf{m} - \mathbf{m}^{(p)} = -\mathbf{B}^{-1}\mathbf{b}$$

#### relate **b** and **B** to g(m)

#### $E(m) = [d-g(m)]^{T}[d-g(m)]$

 $\mathbf{b} = -2\mathbf{G}^{(p)T}[\mathbf{d} - \mathbf{g}(\mathbf{m}^{(p)})]$  and  $\mathbf{B} \approx 2[\mathbf{G}^{(p)T}\mathbf{G}^{(p)}]$ 



#### formula for approximate solution

 $\mathbf{m} - \mathbf{m}^{(p)} = -\mathbf{B}^{-1}\mathbf{b}$ 

 $\mathbf{m} - \mathbf{m}^{(p)} \approx \left[\mathbf{G}^{(p)\mathsf{T}}\mathbf{G}^{(p)}\right]^{-1}\mathbf{G}^{(p)\mathsf{T}}\left[\mathbf{d} - \mathbf{g}(\mathbf{m}^{(p)})\right]$ 



# what do you do if you can't analytically differentiate **g(m)**?

## use finite differences to numerically differentiate

g(m)

or

 $E(\mathbf{m})$ 

#### first derivative

$$\frac{\partial E}{\partial m_{i}}\Big|_{\mathbf{m}^{(p)}} \approx \frac{1}{\Delta m} \left\{ E\left(\mathbf{m} + \Delta \mathbf{m}^{(i)}\right) - E\left(\mathbf{m}\right) \right\}$$

#### first derivative



need to evaluate  $E(\mathbf{m}) M + 1$  times

#### second derivative

$$\frac{\partial^2 E}{\partial m_i \partial m_j} \bigg|_{\mathbf{m}^{(p)}} \approx \frac{1}{(\Delta m)^2} \left\{ E \left( \mathbf{m} + \Delta \mathbf{m}^{(i)} \right) - 2E(\mathbf{m}) + E \left( \mathbf{m} - \Delta \mathbf{m}^{(i)} \right) \right\} \text{ if } i = j$$

$$\frac{\partial^2 E}{\partial m_i \partial m_j} \bigg|_{\mathbf{m}^{(p)}} \approx \frac{1}{4(\Delta m)^2} \left\{ E\left(\mathbf{m} + \Delta \mathbf{m}^{(i)} + \Delta \mathbf{m}^{(j)}\right) - E\left(\mathbf{m} + \Delta \mathbf{m}^{(i)} - \Delta \mathbf{m}^{(j)}\right) - E\left(\mathbf{m} - \Delta \mathbf{m}^{(i)} - \Delta \mathbf{m}^{(j)}\right) - E\left(\mathbf{m} - \Delta \mathbf{m}^{(i)} - \Delta \mathbf{m}^{(j)}\right) + E\left(\mathbf{m} - \Delta \mathbf{m}^{(i)} - \Delta \mathbf{m}^{(j)}\right) \right\} \text{ if } i \neq j$$

need to evaluate  $E(\mathbf{m})$  about  $\frac{1}{2}M^2$  times

#### what can go wrong?

#### convergence to a *local minimum*



#### analytically differentiate sample inverse problem

$$d_i(x_i) = \sin(\omega_0 m_1 x_i) + m_1 m_2$$

$$\mathbf{G}^{(p)} = \begin{bmatrix} \omega_0 x_1 \cos\left(\omega_0 x_1 m_1^{(p)}\right) + m_2^{(p)} & m_1^{(p)} \\ \omega_0 x_2 \cos\left(\omega_0 x_2 m_1^{(p)}\right) + m_2^{(p)} & m_1^{(p)} \\ \vdots & \vdots \\ \omega_0 x_N \cos\left(\omega_0 x_N m_1^{(p)}\right) + m_2^{(p)} & m_1^{(p)} \end{bmatrix}$$



#### often, the convergence is very rapid

#### often, the convergence is very rapid

but sometimes the solution converges to a local minimum and sometimes it even diverges

```
mg = [1, 1]';
G = zeros(N,M);
for k = [1:Niter]
    dg = sin(w0*mg(1)*x) + mg(1)*mg(2);
    dd = dobs - dg;
    Eg=dd'*dd;
    G = zeros(N,2);
    G(:,1) = w0 * x \cdot cos(w0 * mg(1) * x) + mg(2);
    G(:,2) = mg(2) * ones(N,1);
    % least squares solution
    dm = (G'*G) \setminus (G'*dd);
    % update
    mg = mg+dm;
```

#### Part 2

### Newton's Method for an Implicit Theory

Implicit Theory f(d,m)=0

#### with Gaussian prediction error and a priori information about **m**

#### to simplify algebra group **d**, **m** into a vector **x**

 $\mathbf{x} = [\mathbf{d}^{\mathrm{T}}, \mathbf{m}^{\mathrm{T}}]^{\mathrm{T}}$ 

$$< x > = [d^{obsT}, < m >^T]^T$$

$$\begin{bmatrix} \operatorname{cov} \mathbf{x} \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} \operatorname{cov} \mathbf{d} \end{bmatrix} & 0 \\ 0 & \begin{bmatrix} \operatorname{cov} \mathbf{m} \end{bmatrix}_{\mathbf{A}} \end{bmatrix}$$



### represent data and a priori model parameters as a Gaussian $p(\mathbf{x})$

#### f(x)=0 defines a surface in the space of x

#### maximize $p(\mathbf{x})$ on this surface

maximum likelihood point is  $\mathbf{x}^{est}$ 



## can get local maxima if **f(x)** is very non-linear



#### mathematical statement of the problem

minimize  $\Phi = [\mathbf{x} - \langle \mathbf{x} \rangle]^T [\operatorname{cov} \mathbf{x}]^{-1} [\mathbf{x} - \langle \mathbf{x} \rangle]$  subject to  $\mathbf{f}(\mathbf{x}) = 0$ 

# its solution (using Lagrange Multipliers) $[\mathbf{x} - \langle \mathbf{x} \rangle] = [\operatorname{cov} \mathbf{x}]\mathbf{F}^{\mathrm{T}} \{\mathbf{F}[\operatorname{cov} \mathbf{x}]\mathbf{F}^{\mathrm{T}}\}^{-1} \{\mathbf{F}[\mathbf{x} - \langle \mathbf{x} \rangle] - \mathbf{f}(\mathbf{x})\}$ with $F_{ij} = \partial f_i / \partial x_j$

#### mathematical statement of the problem

minimize  $\Phi = [\mathbf{x} - \langle \mathbf{x} \rangle]^T [\operatorname{cov} \mathbf{x}]^{-1} [\mathbf{x} - \langle \mathbf{x} \rangle]$  subject to  $\mathbf{f}(\mathbf{x}) = 0$ 

## its solution (using Lagrange Multipliers) $[\mathbf{x} - \langle \mathbf{x} \rangle] = [\operatorname{cov} \mathbf{x}]\mathbf{F}^{\mathrm{T}} \{\mathbf{F}[\operatorname{cov} \mathbf{x}]\mathbf{F}^{\mathrm{T}}\}^{-1} \{\mathbf{F}[\mathbf{x} - \langle \mathbf{x} \rangle] - \mathbf{f}(\mathbf{x})\}$ reminiscent of minimum length solution

#### mathematical statement of the problem

minimize  $\Phi = [\mathbf{x} - \langle \mathbf{x} \rangle]^T [\operatorname{cov} \mathbf{x}]^{-1} [\mathbf{x} - \langle \mathbf{x} \rangle]$  subject to  $\mathbf{f}(\mathbf{x}) = 0$ 

### its solution (using Lagrange Multipliers) $[x - \langle x \rangle] = [\operatorname{cov} x]F^{T}{F[\operatorname{cov} x]F^{T}}^{-1}{F[x - \langle x \rangle] - f(x)}$

oops! x appears in 3 places

## solution iterate !

$$\mathbf{x}^{(p+1)} = \langle \mathbf{x} \rangle + [\operatorname{cov} \mathbf{x}] \mathbf{F}^{(p)T} \{ \mathbf{F}^{(p)} [\operatorname{cov} \mathbf{x}] \mathbf{F}^{(p)T} \}^{-1} \{ \mathbf{F}^{(p)} [\mathbf{x}^{(p)} - \langle \mathbf{x} \rangle] - \mathbf{f}(\mathbf{x}^{(p)}) \}$$
  
new value for  $\mathbf{x}$  old value for  $\mathbf{x}$  is  $\mathbf{x}^{(p)}$ 

#### special case of an explicit theory f(x) = d-g(m)

$$\mathbf{m}^{(p+1)} = \langle \mathbf{m} \rangle + \mathbf{G}_{(p)}^{-g} \{ \mathbf{d} - \mathbf{g}(\mathbf{m}^{(p)}) + \mathbf{G}^{(p)} [\mathbf{m}^{(p)} - \langle \mathbf{m} \rangle ] \}$$

$$\mathbf{G}_{(p)}^{-g} = \{ [\operatorname{cov} \mathbf{m}]_{A}^{-1} + \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1} \mathbf{G}^{(p)} \}^{-1} \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1}$$
equivalent to solving
$$\begin{bmatrix} [\operatorname{cov} \mathbf{d}]^{-\frac{1}{2}} \mathbf{G}^{(p)} \\ [\operatorname{cov} \mathbf{m}]_{A}^{-\frac{1}{2}} \mathbf{I} \end{bmatrix} \mathbf{m}^{(p+1)} = \begin{bmatrix} [\operatorname{cov} \mathbf{d}]^{-\frac{1}{2}} \{ \mathbf{d} - \mathbf{g}(\mathbf{m}^{(p)}) + \mathbf{G}^{(p)} \mathbf{m}^{(p)} \} \\ [\operatorname{cov} \mathbf{m}]_{A}^{-\frac{1}{2}} \mathbf{I} \end{bmatrix}$$
using simple least squares

#### special case of an explicit theory f(x) = d-g(m)

 $\mathbf{m}^{(p+1)} = \langle \mathbf{m} \rangle + \mathbf{G}_{(p)}^{-g} \{ \mathbf{d} - \mathbf{g}(\mathbf{m}^{(p)}) + \mathbf{G}^{(p)}[\mathbf{m}^{(p)} - \langle \mathbf{m} \rangle] \}$ 

 $\mathbf{G}_{(p)}^{-\mathbf{g}} = \{ [\operatorname{cov} \mathbf{m}]_{A}^{-1} + \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1} \mathbf{G}^{(p)} \}^{-1} \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1}$ weighted least squares generalized inverse with a linearized data kernel

#### special case of an explicit theory f(x) = d-g(m)

$$\mathbf{m}^{(p+1)} = \langle \mathbf{m} \rangle + \mathbf{G}_{(p)}^{-g} \{ \mathbf{d} - \mathbf{g} (\mathbf{m}^{(p)}) + \mathbf{G}^{(p)} [\mathbf{m}^{(p)} - \langle \mathbf{m} \rangle] \}$$

### Newton's Method, but making *E+L* small not just *E* small

 $\mathbf{G}_{(p)}^{-g} = \left\{ [\operatorname{cov} \mathbf{m}]_A^{-1} + \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1} \mathbf{G}^{(p)} \right\}^{-1} \mathbf{G}^{(p)T} [\operatorname{cov} \mathbf{d}]^{-1}$ 

#### Part 3

#### The Gradient Method

# What if you can compute $E(\mathbf{m})$ and $\partial E/\partial m_p$

but you can't compute  $\partial g/\partial m_p$  or  $\partial^2 E/\partial m_p \partial m_q$ 





unit vector pointing towards the minimum  $\mathbf{v} = -\nabla E / |\nabla E|$ 

### so improved solution would be $\mathbf{m}^{(j+1)} = \mathbf{m}^{(j)} + \alpha \mathbf{v}$

if we knew how big to make  $\alpha$ 

#### Armijo's rule provides an acceptance criterion for $\alpha$

$$E(\mathbf{m}^{(k+1)}) \le E(\mathbf{m}^{(k)}) + c\alpha \mathbf{v}^{\mathrm{T}} \nabla E|_{\mathbf{m}^{(k)}} \quad \text{with } c \approx 10^{-4}$$

simple strategy start with a largish  $\alpha$ divide it by 2 whenever it fails Armijo's Rule



```
% error and its gradient at the trial solution
mgo=[1,1]';
ygo = sin(w0*mgo(1)*x) + mgo(1)*mgo(2);
Eqo = (yqo-y) ' * (yqo-y);
dydmo = zeros(N,2);
dydmo(:,1) = w0*x.*cos(w0*mgo(1)*x) + mgo(2);
dydmo(:,2) = mgo(2) * ones(N,1);
dEdmo = 2*dydmo'*(ygo-y);
alpha = 0.05;
c1 = 0.0001;
tau = 0.5;
Niter=500;
for k = [1:Niter]
    v = -dEdmo / sqrt(dEdmo'*dEdmo);
```

```
% backstep
```

```
for kk=[1:10]
    mg = mgo+alpha*v;
    yg = sin(w0*mg(1)*x)+mg(1)*mg(2);
    Eg = (yg-y) ' * (yg-y);
    dydm = zeros(N,2);
    dydm(:,1) = w0*x.*cos(w0*mg(1)*x) + mg(2);
    dydm(:,2) = mg(2) * ones(N,1);
    dEdm = 2*dydm'*(yg-y);
    if( (Eg<=(Ego + c1*alpha*v'*dEdmo)) )</pre>
        break;
    end
    alpha = tau*alpha;
end
```

```
% change in solution
    Dmg = sqrt((mg-mgo)'*(mg-mgo));
    % update
    mgo=mg;
    ygo = yg;
    Ego = Eg;
    dydmo = dydm;
    dEdmo = dEdm;
    if (Dmg < 1.0e-6)
        break;
    end
end
```

# often, the convergence is reasonably rapid

# often, the convergence is reasonably rapid

#### exception

when the minimum is in along a long shallow valley