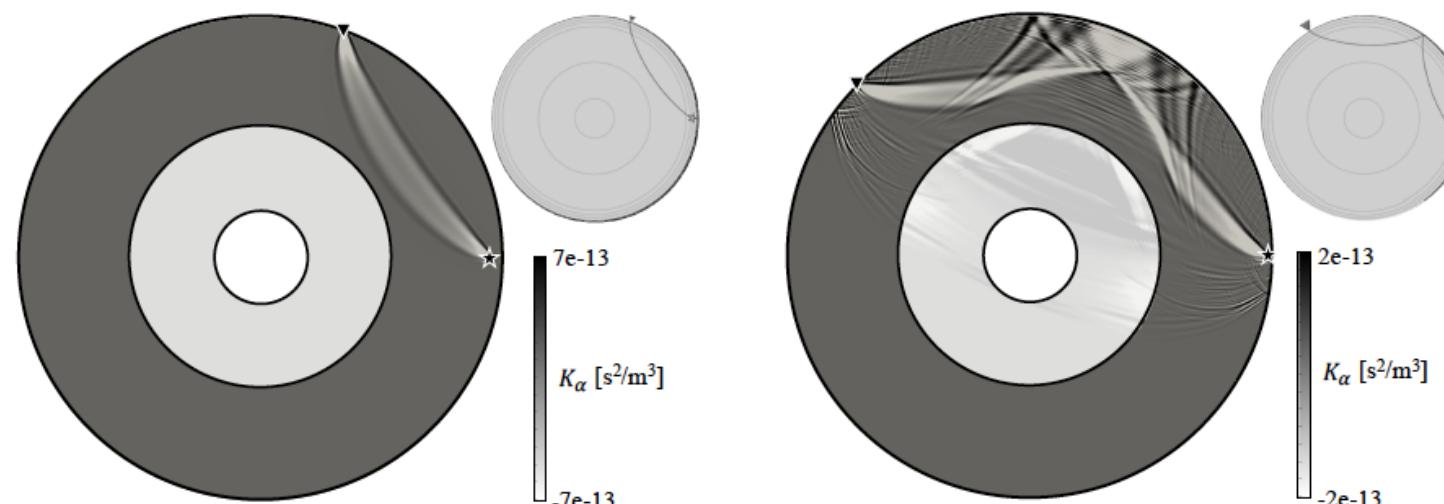


## Lecture 2

# Non-linear optimisation and adjoint methods

Andreas Fichtner

ETH Zurich - Seismology and Wave Physics



# MOTIVATION

- More often than not, the forward model is not linear.
- The beautiful machinery of linear inverse theory is not applicable.
- We still would like to find a model that explains observations acceptably well.

# GOALS OF THIS LECTURE

- Introduce basic concepts of non-linear optimisation.
- Cover the most important classes of non-linear, gradient-based optimisation methods.
- Introduce adjoint methods for the efficient computation of misfit gradients.

# OUTLINE

## PART I: Basic non-linear optimisation concepts

1. Iterative linearisation
2. Local optimisation, global and local minima

## PART II: Non-linear optimisation methods

1. General gradient descent methods
2. The method of steepest descent
3. Newton's method
4. Other methods: Conjugate gradients, BFGS, L-BFGS, ...

## PART III: Non-linear optimisation Jupyter notebook

## PART IV: Adjoint methods

1. Motivation
2. Discrete adjoints
3. Continuous adjoints
4. Sensitivity kernels and kernel gallery

## PART V: Adjoint method Jupyter notebook

All of this will closely follow the *Fundamentals of Inverse Theory* lecture notes.

## PART I

Basic non-linear optimisation concepts

## 1. Iterative linearisation

- Frequently, the forward problem  $\mathbf{d}=\mathbf{G}(\mathbf{m})$  is not linear.
- The least-squares method develop for linear forward problems is not applicable.
  - No explicit solution for the posterior mean model and posterior covariance.
  - No resolution matrix.
  - ...

- Frequently, the forward problem  $\mathbf{d}=\mathbf{G}(\mathbf{m})$  is not linear.
- The least-squares method develop for linear forward problems is not applicable.
  - No explicit solution for the posterior mean model and posterior covariance.
  - No resolution matrix.
  - ...
- The seemingly most straightforward solution: Pick a prior model  $\mathbf{m}_0$  and linearise:

$$\mathbf{G}(\mathbf{m}) \doteq \mathbf{G}(\mathbf{m}_0) + \mathbf{J}_0(\mathbf{m} - \mathbf{m}_0) \quad J_{0,ij} = \left. \frac{\partial G_i}{\partial m_j} \right|_{\mathbf{m}=\mathbf{m}_0} = \text{Jacobian at } \mathbf{m}_0$$

- Now this is a linear problem! We can apply the linear least-squares machinery to obtain a ‘solution’  $\mathbf{m}_1$ .

- Frequently, the forward problem  $\mathbf{d}=\mathbf{G}(\mathbf{m})$  is not linear.
- The least-squares method develop for linear forward problems is not applicable.
  - No explicit solution for the posterior mean model and posterior covariance.
  - No resolution matrix.
  - ...
- The seemingly most straightforward solution: Pick a prior model  $\mathbf{m}_0$  and linearise:

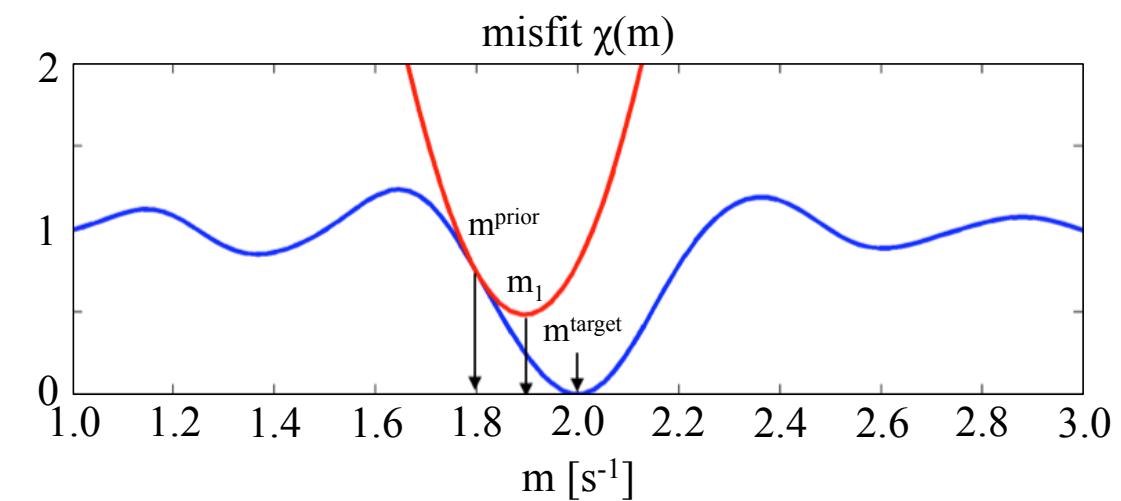
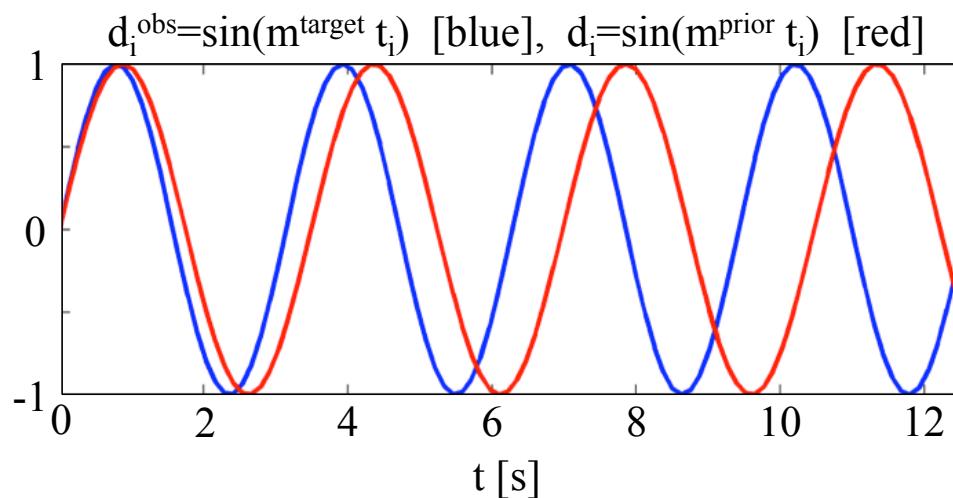
$$\mathbf{G}(\mathbf{m}) \doteq \mathbf{G}(\mathbf{m}_0) + \mathbf{J}_0(\mathbf{m} - \mathbf{m}_0) \quad J_{0,ij} = \left. \frac{\partial G_i}{\partial m_j} \right|_{\mathbf{m}=\mathbf{m}_0} = \text{Jacobian at } \mathbf{m}_0$$

- Now this is a linear problem! We can apply the linear least-squares machinery to obtain a ‘solution’  $\mathbf{m}_1$ .
- This can be **iterated**:
  - Linearise forward problem around  $\mathbf{m}_1$  using new Jacobian  $\mathbf{J}_1$ .
  - Compute new least-squares solution  $\mathbf{m}_2$ .
  - Repeat ...

# ITERATIVE LINEARISATION - EXAMPLE

- A toy problem:

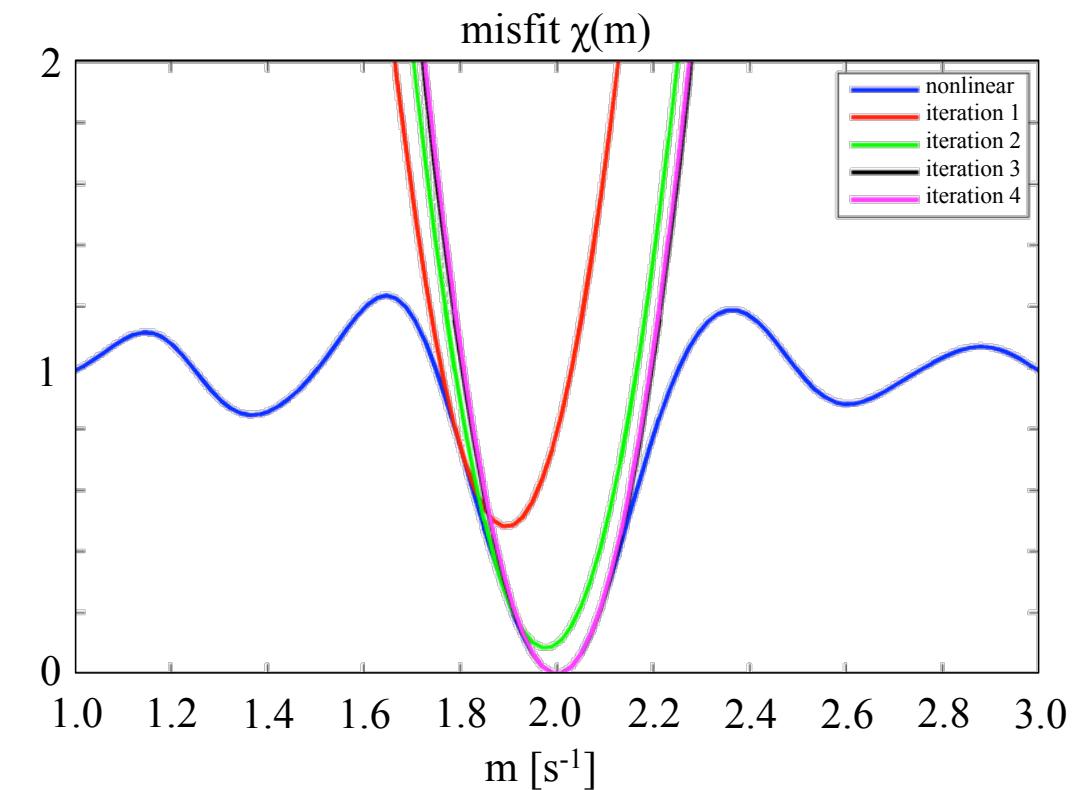
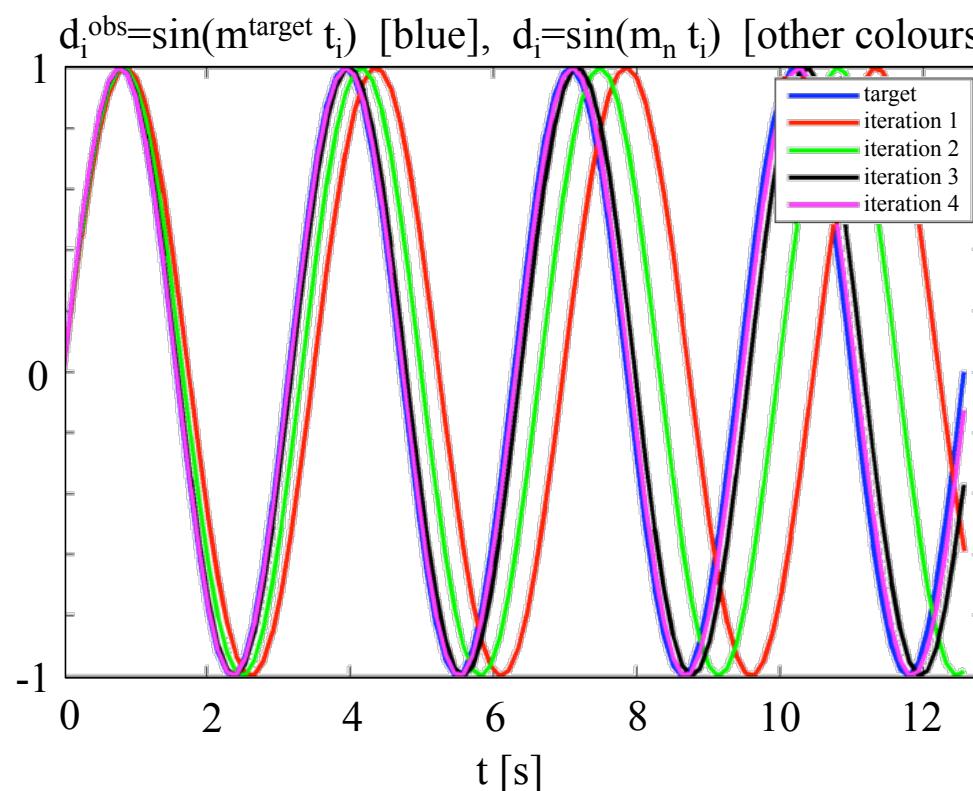
$$d_i = \sin(mt_i) \quad \chi(m) = \frac{1}{2} \sum_{i=1}^N \left[ d_i^{\text{obs}} - \sin(mt_i) \right]^2$$



# ITERATIVE LINEARISATION - EXAMPLE

- A toy problem:

$$d_i = \sin(mt_i) \quad \chi(m) = \frac{1}{2} \sum_{i=1}^N \left[ d_i^{\text{obs}} - \sin(mt_i) \right]^2$$

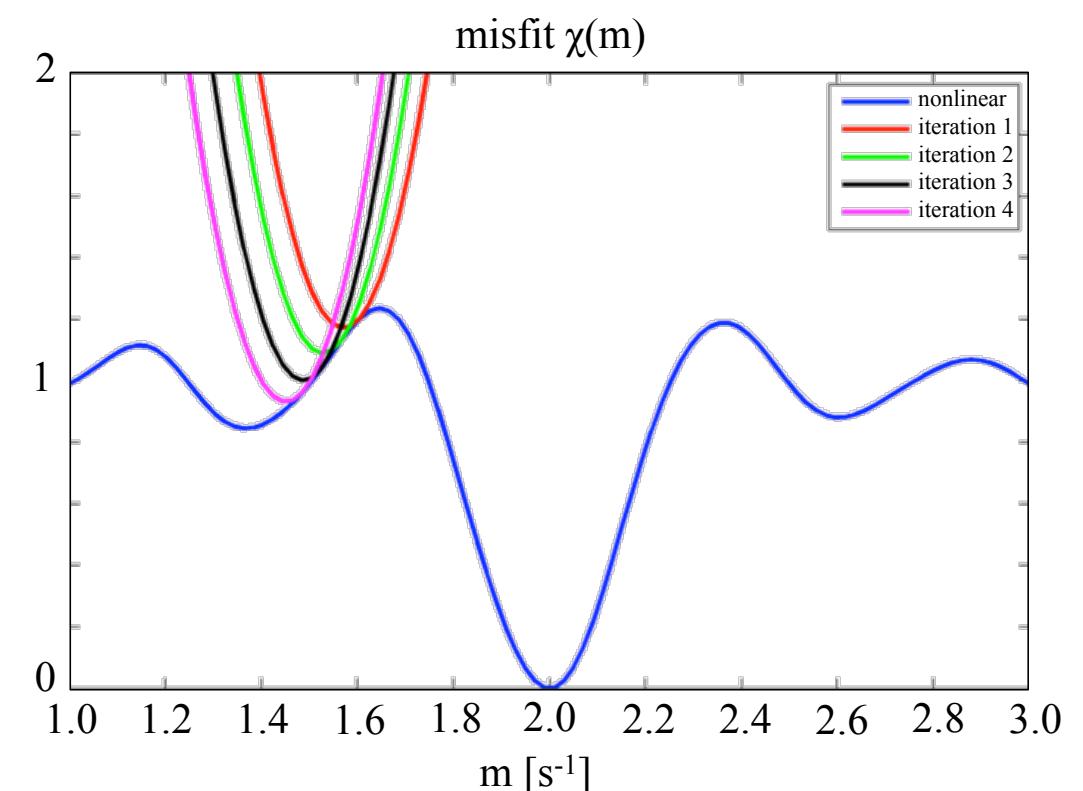
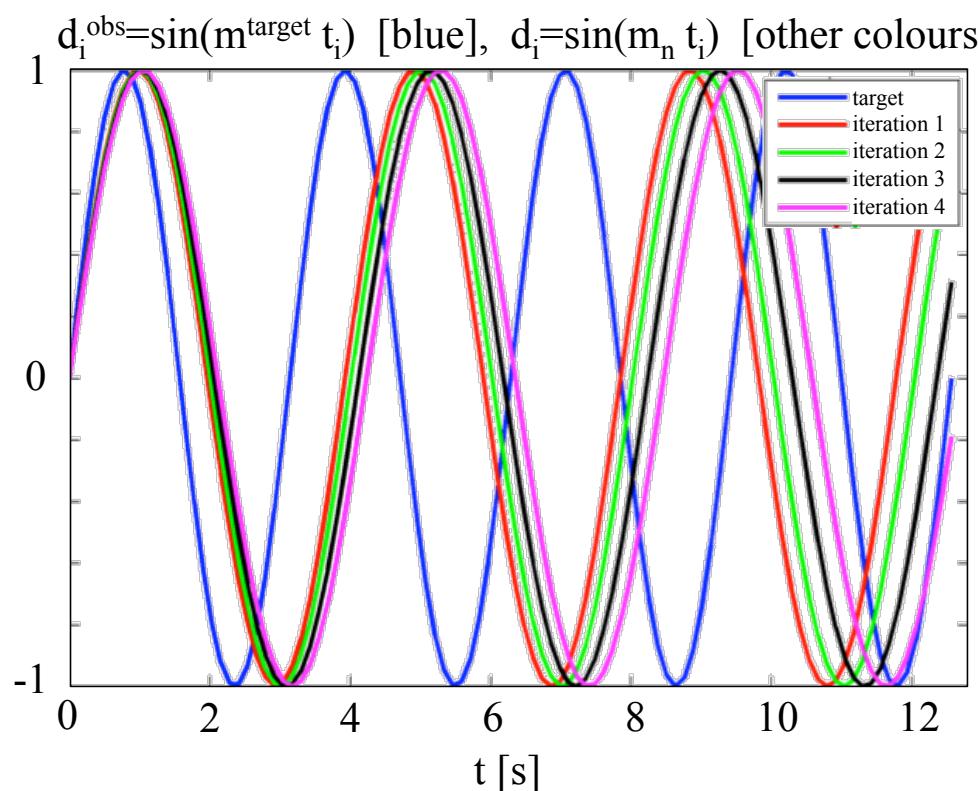


## 2. Local optimisation, global and local minima

# ITERATIVE LINEARISATION – EFFECT OF INITIAL MODEL CHOICE

- A toy problem:

$$d_i = \sin(mt_i) \quad \chi(m) = \frac{1}{2} \sum_{i=1}^N \left[ d_i^{\text{obs}} - \sin(mt_i) \right]^2$$



- A toy problem:

$$d_i = \sin(mt_i) \quad \chi(m) = \frac{1}{2} \sum_{i=1}^N \left[ d_i^{\text{obs}} - \sin(mt_i) \right]^2$$

- Iterative linearisation [and all methods we will see in this lecture] is a local method.
- It may get stuck in a local minimum.
- The ‘result’ may depend strongly on the choice of the initial model.
- In practice, it may be very hard to check if the minimum is local or global.
  - Analysis of data fit
  - Plausibility arguments
  - Check against independent data

## **PART II**

Non-linear optimisation methods

## 1. General descent methods

## THE GENERAL SCHEME

- Our goal: Find a model  $\mathbf{m}$  such that a suitably defined misfit  $\chi$  is minimal.
- The minimisation proceeds iteratively:
  1. Start from initial model  $\mathbf{m}_0$
  2. Update according to  $\mathbf{m}_{i+1} = \mathbf{m}_i + \gamma_i \mathbf{h}_i$  with  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$   
step length  

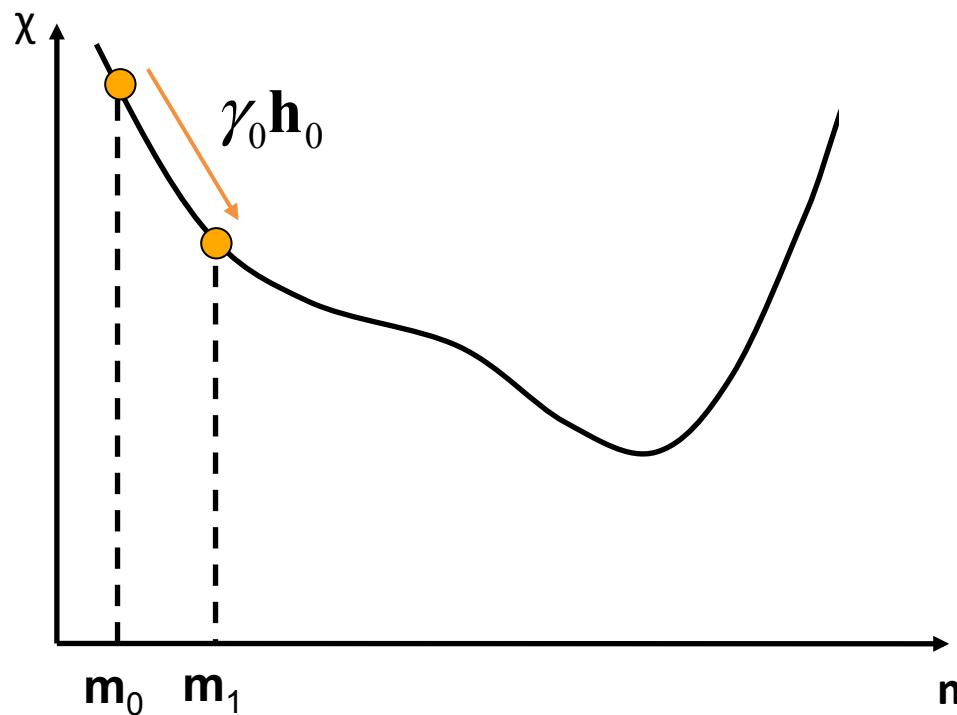
## THE GENERAL SCHEME

- Our goal: Find a model  $\mathbf{m}$  such that a suitably defined misfit  $\chi$  is minimal.
- The minimisation proceeds iteratively:

1. Start from initial model  $\mathbf{m}_0$

2. Update according to  $\mathbf{m}_{i+1} = \mathbf{m}_i + \gamma_i \mathbf{h}_i$  with  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$

step length      descent direction



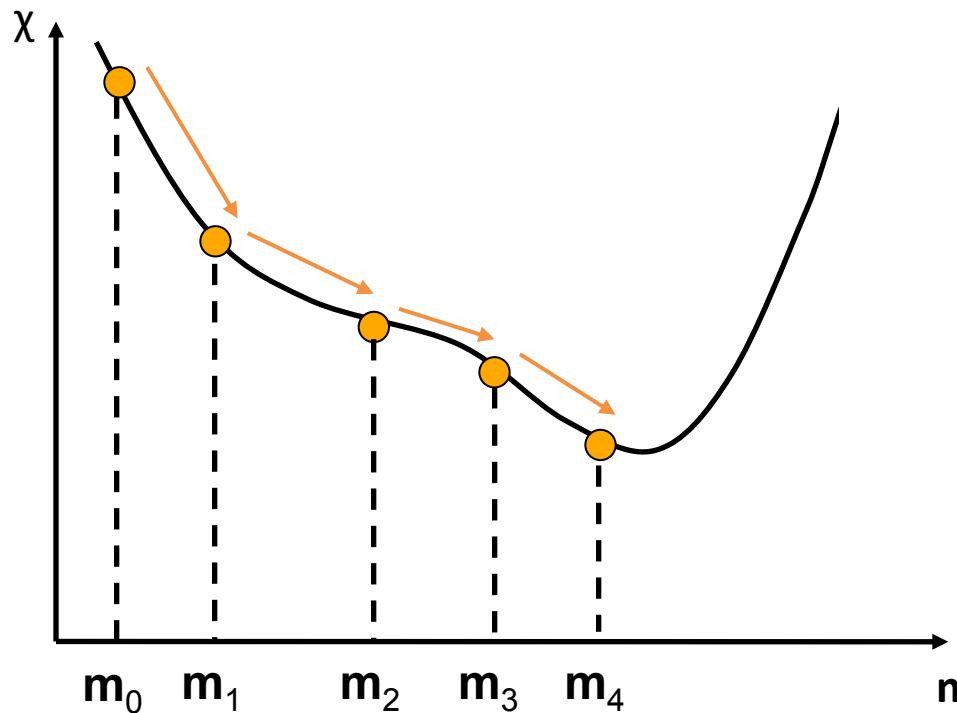
## THE GENERAL SCHEME

- Our goal: Find a model  $\mathbf{m}$  such that a suitably defined misfit  $\chi$  is minimal.
- The minimisation proceeds iteratively:

1. Start from initial model  $\mathbf{m}_0$

2. Update according to  $\mathbf{m}_{i+1} = \mathbf{m}_i + \gamma_i \mathbf{h}_i$  with  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$

step length      descent direction



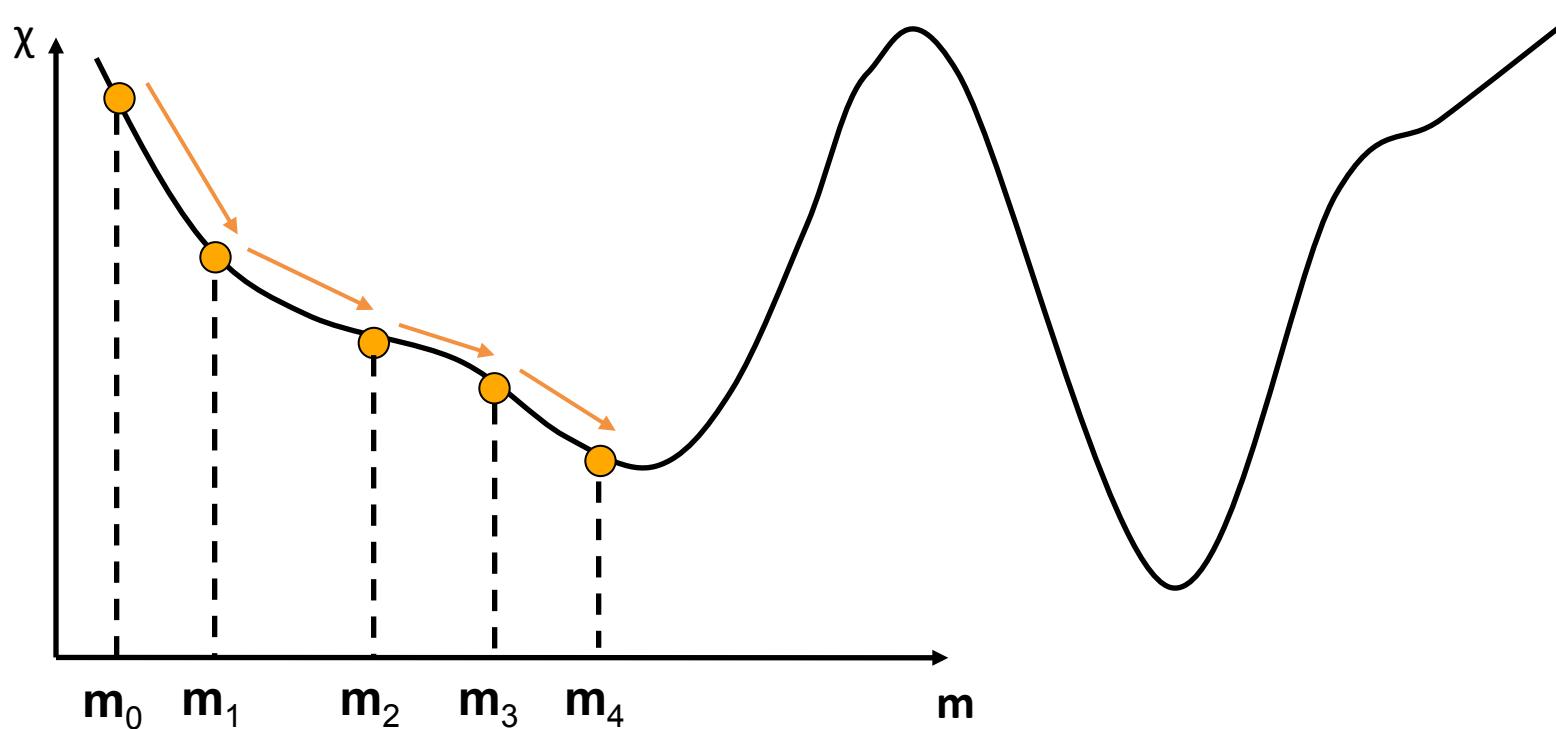
## THE GENERAL SCHEME

- Our goal: Find a model  $\mathbf{m}$  such that a suitably defined misfit  $\chi$  is minimal.
- The minimisation proceeds iteratively:

1. Start from initial model  $\mathbf{m}_0$

2. Update according to  $\mathbf{m}_{i+1} = \mathbf{m}_i + \gamma_i \mathbf{h}_i$  with  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$

step length      descent direction



- Our goal: Find a model  $\mathbf{m}$  such that a suitably defined misfit  $\chi$  is minimal.
- The minimisation proceeds iteratively:

1. Start from initial model  $\mathbf{m}_0$

2. Update according to  $\mathbf{m}_{i+1} = \mathbf{m}_i + \gamma_i \mathbf{h}_i$  with  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$

step length  descent direction

- A possible descent direction is always:  $\mathbf{h}_0 = -\nabla \chi(\mathbf{m}_0)$
- Other possible directions are  $\mathbf{h}_0 = -\mathbf{A} \nabla \chi(\mathbf{m}_0)$  for any positive definite matrix  $\mathbf{A}$  [which may change with iteration].
- The members of the descent method family are defined by the choice of  $\mathbf{A}$ .

## 2. The method of steepest descent

- $\mathbf{A}=\mathbf{I} \rightarrow \mathbf{h}$  is direction of steepest descent

### Method of steepest descent

1. Choose an initial model,  $\mathbf{m}_0$ , and set  $i = 0$ .
2. Compute the gradient for the current model,  $\nabla\chi(\mathbf{m}_i)$ .
3. Update  $\mathbf{m}_i$  according to

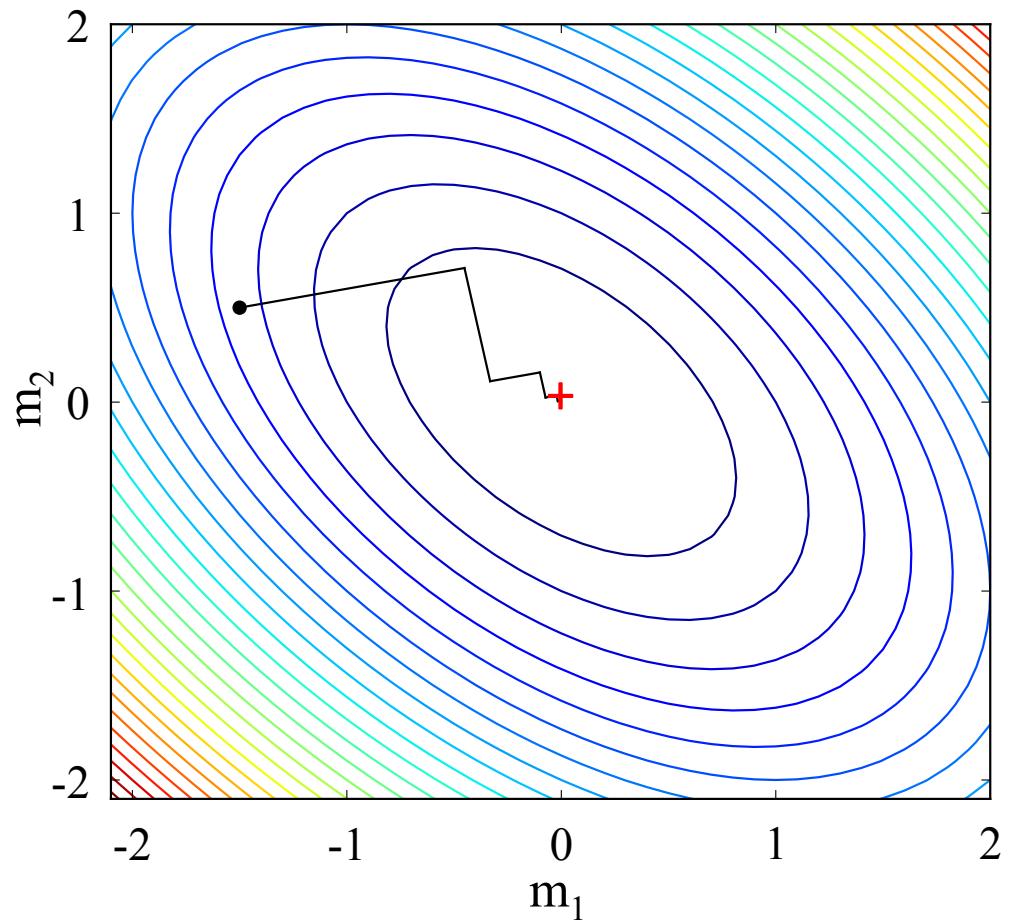
$$\mathbf{m}_{i+1} = \mathbf{m}_i - \gamma_i \nabla\chi(\mathbf{m}_i),$$

with a suitable step length  $\gamma_i$  that ensures  $\chi(\mathbf{m}_{i+1}) < \chi(\mathbf{m}_i)$ .

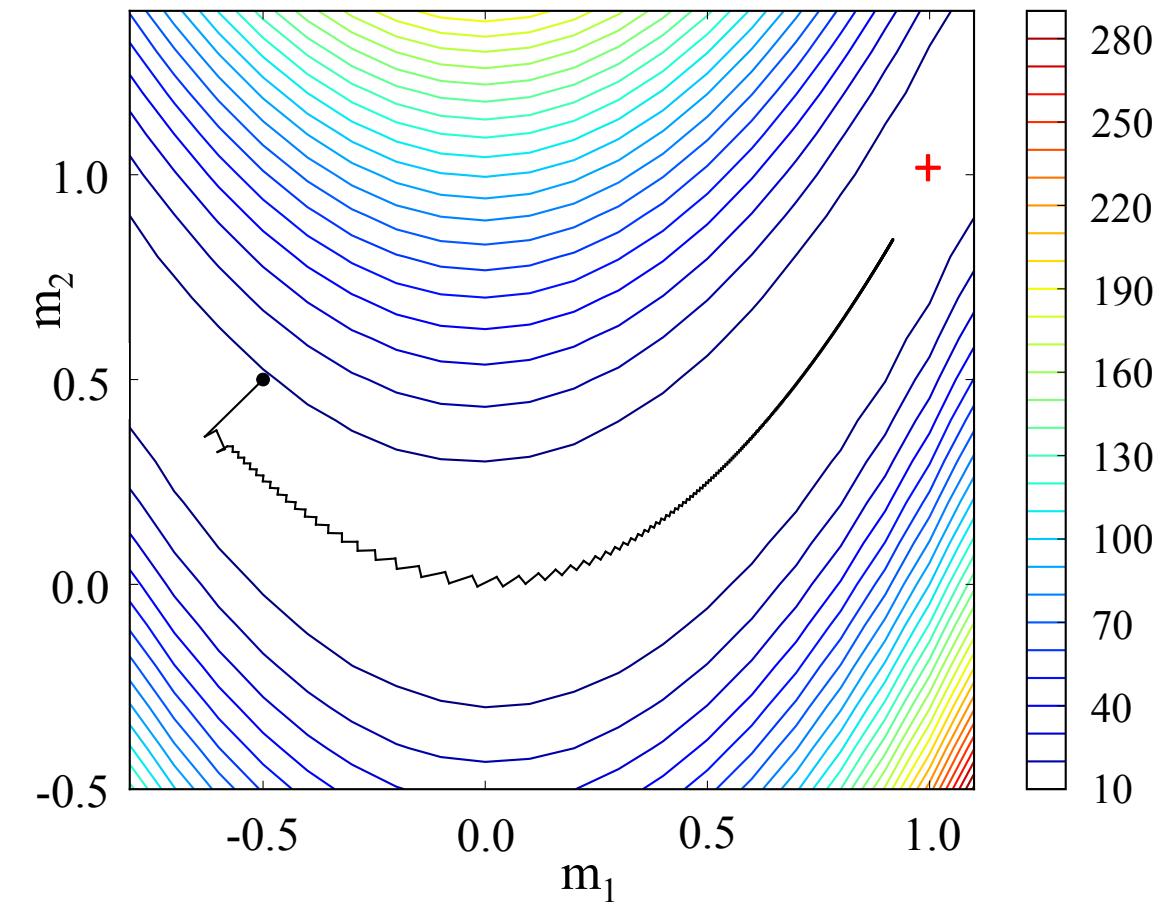
4. Set  $i \rightarrow i + 1$ , go back to (2) and repeat until the data are explained sufficiently well.

- $\mathbf{A}=\mathbf{I} \rightarrow \mathbf{h}$  is direction of steepest descent

a) quadratic function



b) Rosenbrock function



## STEEPEST DESCENT – PROS & CONS

- $\mathbf{A}=\mathbf{I}$  →  $\mathbf{h}$  is direction of steepest descent
- + Very easy to implement.
- + Method of choice for testing purposes.
- Potentially very slow convergence.

### 3. Newton's method

## NEWTON'S METHOD – THE ALGORITHM

- Newton's method is motivated by a Taylor expansion:

$$\mathbf{0} = \nabla \chi(\tilde{\mathbf{m}}) \approx \nabla \chi(\mathbf{m}) + \mathbf{H}(\mathbf{m})(\tilde{\mathbf{m}} - \mathbf{m}) \quad \tilde{\mathbf{m}} \approx \mathbf{m} - \mathbf{H}^{-1}(\mathbf{m}) \nabla \chi(\mathbf{m})$$

- $\mathbf{A}=\mathbf{H}^{-1}$ , where  $\mathbf{H}$  is the **Hessian** [matrix of second derivatives].

- Newton's method is motivated by a Taylor expansion:

$$\mathbf{0} = \nabla \chi(\tilde{\mathbf{m}}) \approx \nabla \chi(\mathbf{m}) + \mathbf{H}(\mathbf{m})(\tilde{\mathbf{m}} - \mathbf{m}) \quad \tilde{\mathbf{m}} \approx \mathbf{m} - \mathbf{H}^{-1}(\mathbf{m}) \nabla \chi(\mathbf{m})$$

- $\mathbf{A}=\mathbf{H}^{-1}$ , where  $\mathbf{H}$  is the **Hessian** [matrix of second derivatives].

### Newton's method

1. Choose an initial model,  $\mathbf{m}_0$ , and set  $i = 0$ .
2. Compute the gradient for the current model,  $\nabla \chi(\mathbf{m}_i)$ .
3. Determine the descent direction,  $\mathbf{h}_i$ , as the solution of

$$\mathbf{H}(\mathbf{m}_i) \mathbf{h}_i = -\nabla \chi(\mathbf{m}_i).$$

4. Update  $\mathbf{m}_i$  according to

$$\mathbf{m}_{i+1} = \mathbf{m}_i + \mathbf{h}_i.$$

5. Set  $i \rightarrow i + 1$ , go back to (2) and repeat as often as needed.

## NEWTON'S METHOD – EXAMPLES

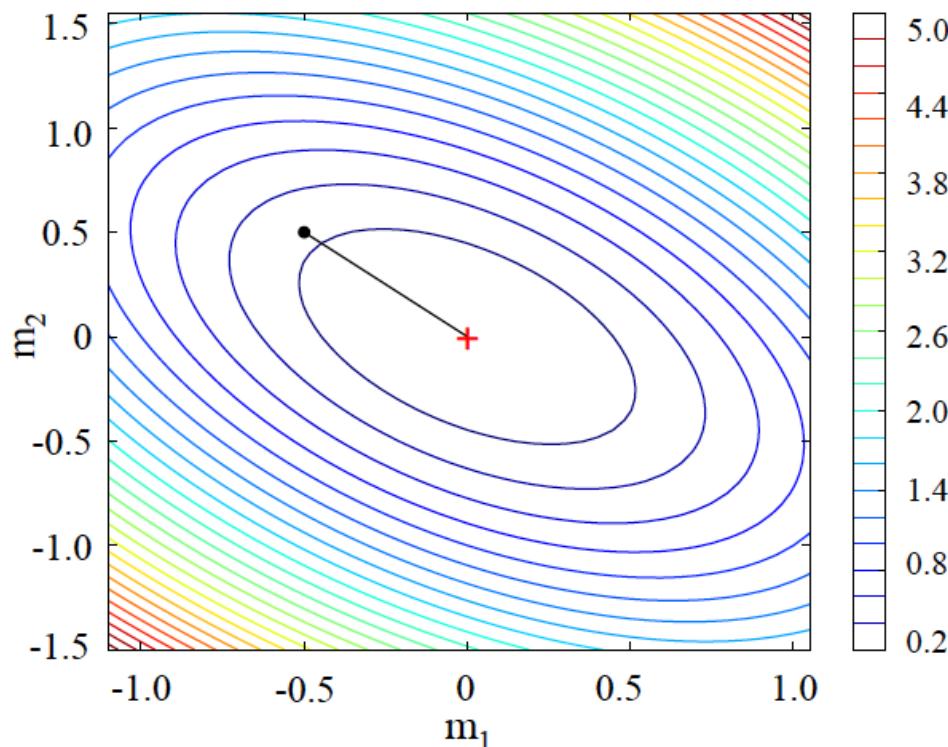
- Newton's method is motivated by a Taylor expansion:

$$\mathbf{0} = \nabla \chi(\tilde{\mathbf{m}}) \approx \nabla \chi(\mathbf{m}) + \mathbf{H}(\mathbf{m})(\tilde{\mathbf{m}} - \mathbf{m})$$

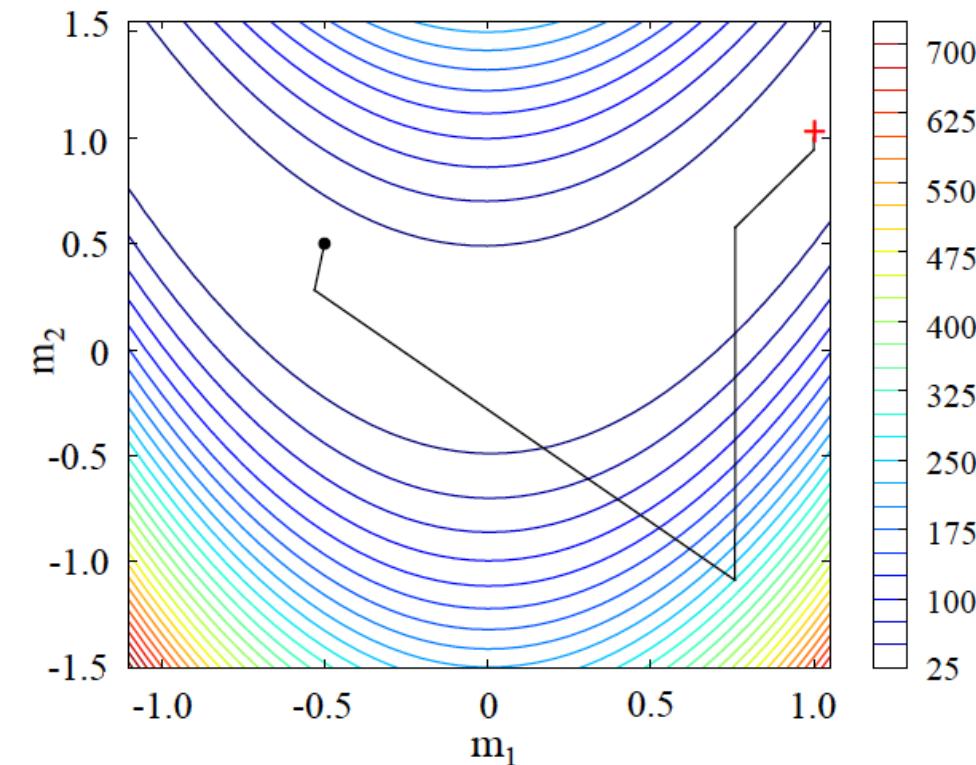
$$\tilde{\mathbf{m}} \approx \mathbf{m} - \mathbf{H}^{-1}(\mathbf{m}) \nabla \chi(\mathbf{m})$$

- $\mathbf{A}=\mathbf{H}^{-1}$ , where  $\mathbf{H}$  is the **Hessian** [matrix of second derivatives].

a) quadratic function



b) Rosenbrock function



- Newton's method is motivated by a Taylor expansion:

$$\mathbf{0} = \nabla \chi(\tilde{\mathbf{m}}) \approx \nabla \chi(\mathbf{m}) + \mathbf{H}(\mathbf{m})(\tilde{\mathbf{m}} - \mathbf{m}) \quad \tilde{\mathbf{m}} \approx \mathbf{m} - \mathbf{H}^{-1}(\mathbf{m}) \nabla \chi(\mathbf{m})$$

- $\mathbf{A}=\mathbf{H}^{-1}$ , where  $\mathbf{H}$  is the **Hessian** [matrix of second derivatives].

+ Converges rapidly. In 1 iteration, when this misfit is quadratic.

- Requires computation of second derivatives.  
- Requires inversion of a big matrix that may not be invertible.

4. Other methods: Conjugate gradients, BFGS, L-BFGS, ...

## THE SPECTRUM OF DESCENT METHODS

- Steepest descent and Newton's method: end members in a broad spectrum of gradient-based descent methods



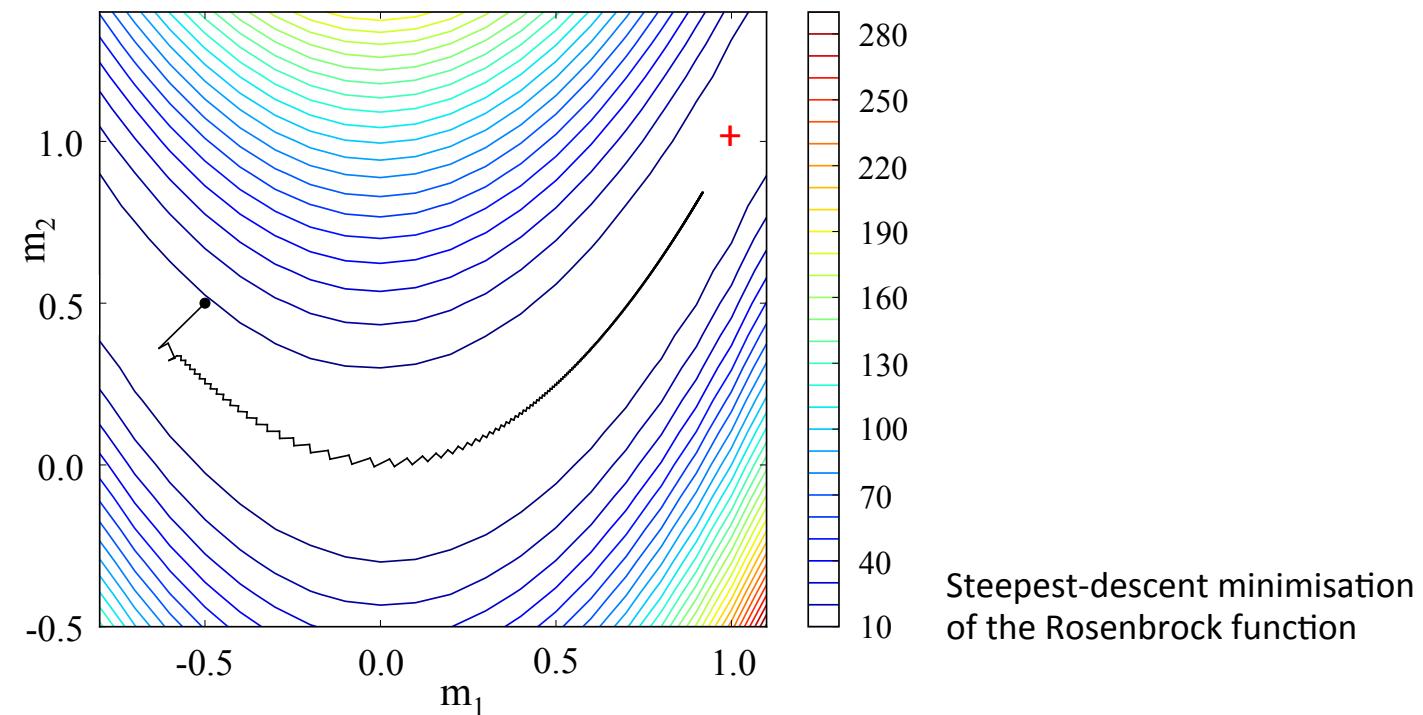
## THE SPECTRUM OF DESCENT METHODS

- Steepest descent and Newton's method: end members in a broad spectrum of gradient-based descent methods



# CONJUGATE GRADIENTS

- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.



- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.
- Basic idea of conjugate gradients: March in directions that are **mutually orthogonal** in the sense of  $\mathbf{h}_i \mathbf{H} \mathbf{h}_{i+1} = 0$ .

- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.
- Basic idea of conjugate gradients: March in directions that are **mutually orthogonal** in the sense of  $\mathbf{h}_i \mathbf{H} \mathbf{h}_{i+1} = 0$ .
- The conjugate-gradient method constructs **a sequence of orthogonal descent directions**,  $\mathbf{h}_0, \mathbf{h}_1, \dots$ , **automatically**.
- There are different ways of doing that, and each of these ways has its own name:
  - Fletcher-Reves method
  - Hestenes-Stiefel method
  - Polak-Ribiere method

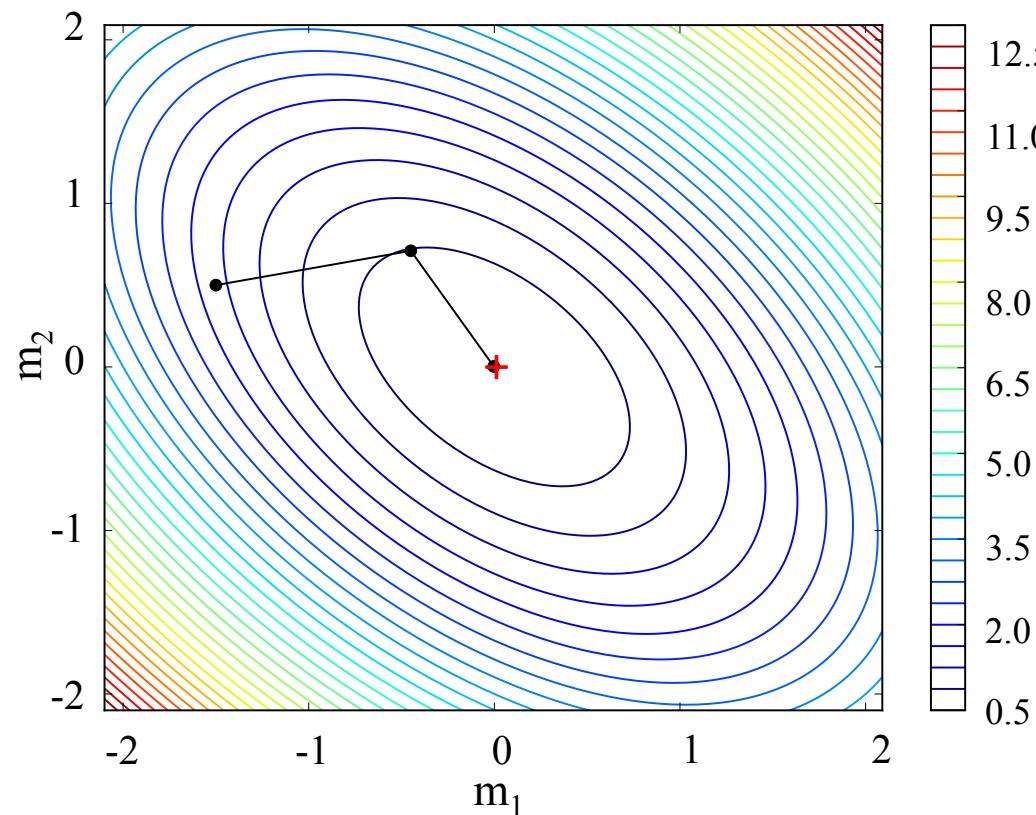
- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.
- Basic idea of conjugate gradients: March in directions that are **mutually orthogonal** in the sense of  $\mathbf{h}_i \mathbf{H} \mathbf{h}_{i+1} = 0$ .
- The conjugate-gradient method constructs **a sequence of orthogonal descent directions**,  $\mathbf{h}_0, \mathbf{h}_1, \dots$ , **automatically**.
- There are different ways of doing that, and each of these ways has its own name:
  - Fletcher-Reves method
  - Hestenes-Stiefel method
  - Polak-Ribiere method
- The conjugate-gradient method **for exactly quadratic misfit functions**:
  - **Converges in  $n$  steps** at most because:
  - **$\mathbf{H}$  is constant** and there are **only  $n$**  mutually orthogonal directions.

- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.
- Basic idea of conjugate gradients: March in directions that are **mutually orthogonal** in the sense of  $\mathbf{h}_i \mathbf{H} \mathbf{h}_{i+1} = 0$ .
- The conjugate-gradient method constructs **a sequence of orthogonal descent directions**,  $\mathbf{h}_0, \mathbf{h}_1, \dots$ , **automatically**.
- There are different ways of doing that, and each of these ways has its own name:
  - Fletcher-Reves method
  - Hestenes-Stiefel method
  - Polak-Ribiere method
- The conjugate-gradient method **for exactly quadratic misfit functions**:
  - **Converges in  $n$  steps** at most because:
  - $\mathbf{H}$  is constant and there are only  $n$  mutually orthogonal directions.
- The conjugate-gradient method **for general misfit functions**:
  - Is basically a leap of faith.
  - Tends to converge more quickly than steepest descent for many applications.

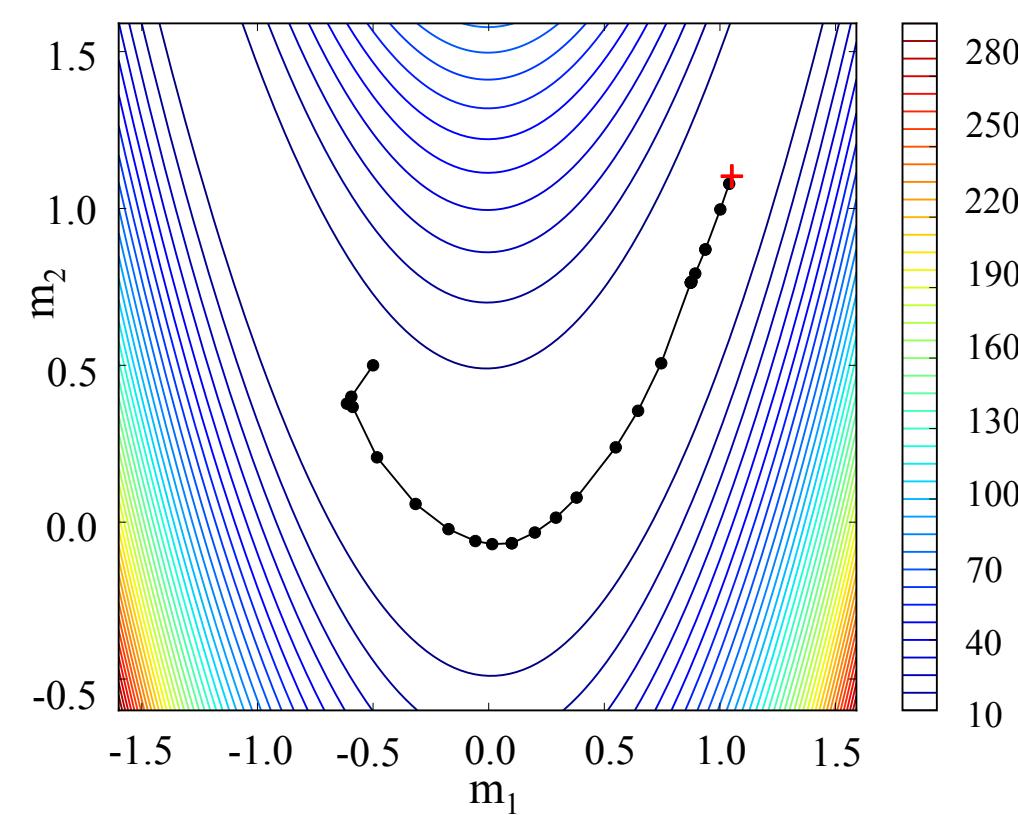
## CONJUGATE GRADIENTS

- Steepest descent tends to go again and again in a similar direction. → Makes it pretty inefficient.
- Basic idea of conjugate gradients: March in directions that are **mutually orthogonal** in the sense of  $\mathbf{h}_i \mathbf{H} \mathbf{h}_{i+1} = 0$ .

a) quadratic function



b) Rosenbrock function





- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

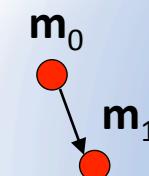
$0^{\text{th}}$  approximation of  $\mathbf{H}^{-1} = \mathbf{I}$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

1<sup>st</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_1$

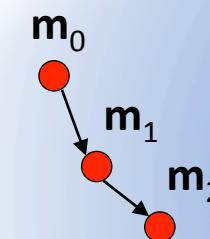
based on  $\mathbf{m}_0, \mathbf{m}_1$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

2<sup>nd</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_2$

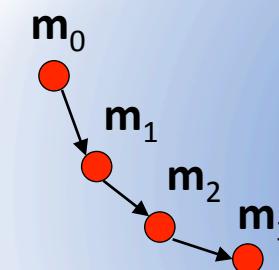
based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

3<sup>rd</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_3$

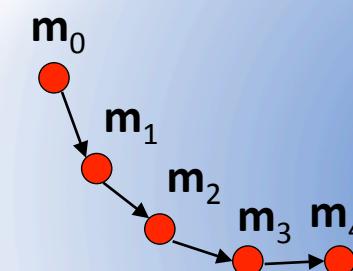
based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

4<sup>th</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_4$

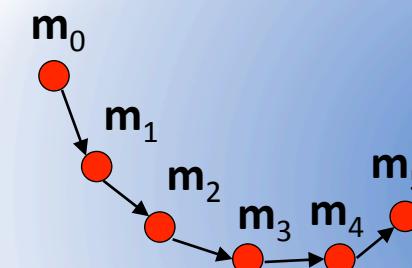
based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}_4$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

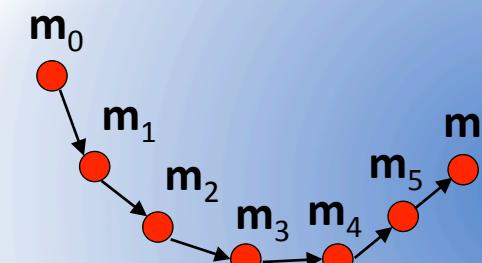
5<sup>th</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_5$

based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}_4, \mathbf{m}_5$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

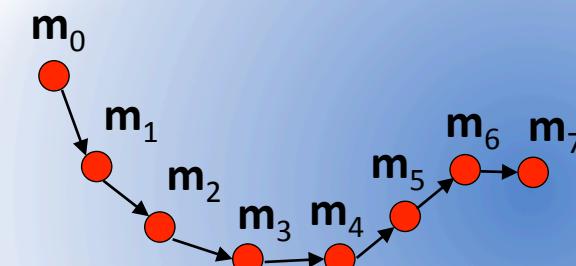
6<sup>th</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_6$   
based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}_4, \mathbf{m}_5, \mathbf{m}_6$



- The Newton method **converges quickly** because it uses **second-derivative** information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.

7<sup>th</sup> approximation of  $\mathbf{H}^{-1} = \mathbf{A}_7$

based on  $\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}_4, \mathbf{m}_5, \mathbf{m}_6, \mathbf{m}_7$



- The Newton method converges quickly because it uses second-derivative information in the form of  $\mathbf{H}^{-1}$ .
- Second derivatives are really difficult to compute! And  $\mathbf{H}$  is difficult to invert!
- Basic idea of the BFGS method: Construct an approximation of  $\mathbf{H}^{-1}$  as we march through model space.
- For exactly quadratic misfit functions:
  - $\mathbf{A}_n = \mathbf{H}^{-1}$  [after  $n$  iterations we get the exact inverse Hessian]
  - Descent directions are **mutually orthogonal**, as in conjugate-gradient methods → **convergence after  $n$  iterations**.
- For general misfit functions: leap of faith but empirically very efficient in many cases.
- Behaviour for test functions: See Jupyter notebook!



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

$0^{\text{th}}$  approximation of  $\mathbf{H}^{-1} = \mathbf{I}$

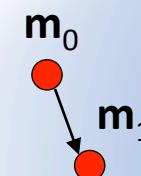
$\mathbf{m}_0$

A 2D plot showing concentric blue circles centered at a red dot labeled  $\mathbf{m}_0$ . The background is white, and the circles are light blue.

- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

1<sup>st</sup> approximation of  $H^{-1} = A_1$

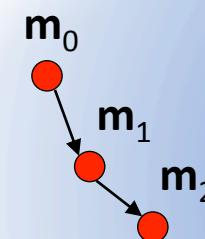
based on  $m_0, m_1$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

2<sup>nd</sup> approximation of  $H^{-1} = A_2$

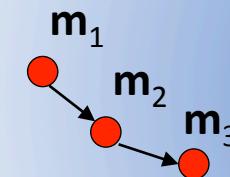
based on  $m_0, m_1, m_2$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

3<sup>rd</sup> approximation of  $H^{-1} = A_3$

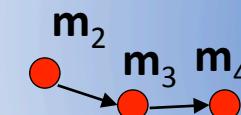
based on  $m_1, m_2, m_3$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

4<sup>th</sup> approximation of  $H^{-1} = A_4$

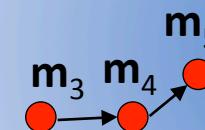
based on  $m_2, m_3, m_4$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

5<sup>th</sup> approximation of  $H^{-1} = A_5$

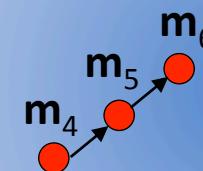
based on  $m_3, m_4, m_5$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

6<sup>th</sup> approximation of  $H^{-1} = A_6$

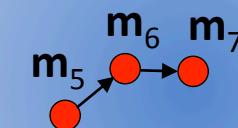
based on  $m_4, m_5, m_6$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!

7<sup>th</sup> approximation of  $H^{-1} = A_7$

based on  $m_5, m_6, m_7$



- The BFGS method requires a lot of storage [matrix  $A_i$ ].
- Basic idea of the L-BFGS method: Just do not use all of the models along the path!
- More difficult to implement than the other methods.
- Empirically hard to beat in the large-scale inverse problems that we deal with!
- Behavior for test functions: See Jupyter Notebook!

# **PART IV**

## Adjoint methods

## 1. Motivation

## So, WHERE IS THE PROBLEM?

- The full gradient – with all its components - is needed in each iteration.
- The most straightforward approach: approximate the gradient by finite-differences:

$$\frac{\partial \chi(m)}{\partial m_k} \approx \frac{\chi(..., m_k + \delta m, ...) - \chi(..., m_k, ...)}{\delta m}$$

- Example with 500,000 model parameters:

500,001 forward simulations

- × 0.5 h per simulation
- × 126 compute cores
- × 50 sources (earthquakes)
- × 50 conjugate gradient iterations

**78e<sup>9</sup> cpu hours ≈ 8,900,000 cpu years**

## 2. Discrete adjoints

## DISCRETE ADJOINT METHOD SUMMARY

Regular wave equation

$$\underline{\underline{L}} \underline{u} = \underline{f}$$

Adjoint wave equation

$$\underline{\underline{L}}^T \underline{v} = -\nabla \chi$$

Gradient equation

$$\frac{\partial \chi}{\partial m_i} = \underline{v}^T \frac{\partial \underline{\underline{L}}}{\partial m_i} \underline{u}$$

Regular wave equation

$$\underline{\underline{L}} \underline{u} = \underline{f}$$

Adjoint wave equation

$$\underline{\underline{L}}^T \underline{v} = -\nabla \chi$$

Gradient equation

$$\frac{\partial \chi}{\partial m_i} = \underline{v}^T \frac{\partial \underline{\underline{L}}}{\partial m_i} \underline{u}$$

**Adjoint recipe**

1. Solve forward problem [regular wave equation] to obtain  $\underline{u}$ .
2. Evaluate misfit  $\chi$ .
3. Compute adjoint source,  $-\nabla \chi$ .
4. Solve adjoint equation to obtain adjoint field  $\underline{v}$ .
5. Plug  $\underline{u}$  and  $\underline{v}$  into the gradient equation.

Regular wave equation

$$\underline{\underline{L}} \underline{u} = \underline{f}$$

Adjoint wave equation

$$\underline{\underline{L}}^T \underline{v} = -\nabla \chi$$

Gradient equation

$$\frac{\partial \chi}{\partial m_i} = \underline{v}^T \frac{\partial \underline{\underline{L}}}{\partial m_i} \underline{u}$$

**Adjoint recipe**

1. Solve forward problem [regular wave equation] to obtain  $\underline{u}$ .
2. Evaluate misfit  $\chi$ .
3. Compute adjoint source,  $-\nabla \chi$ .
4. Solve adjoint equation to obtain adjoint field  $\underline{v}$ .
5. Plug  $\underline{u}$  and  $\underline{v}$  into the gradient equation.

**Comments**

1. No need to explicitly compute the derivative of the wavefield  $\underline{u}$  [by construction].
2. Gradient is entirely determined by the definition of the misfit [adjoint source is the only thing that explicitly depends on the misfit].
3. Computation of gradient requires storage of forward wavefield  $\underline{u}$ .

### 3. Continuous adjoints

Discrete case [frequency domain]

$$\underline{\underline{L}}\underline{u} = (-\omega^2 \underline{\underline{M}} + \underline{\underline{K}})\underline{u}$$

$$\nabla \chi = \underline{\underline{v}}^T \nabla \underline{\underline{L}} \underline{u}$$

Continuous case [time domain]

$$L(\underline{u}) = \rho \ddot{\underline{u}} - \nabla \cdot (\mathbf{C} \cdot \nabla \underline{u}) = \underline{f}$$

$$\nabla \chi' = \int \underline{\underline{v}}^T \nabla L(\underline{u}) dt$$

- The same formal derivation from the discrete case can be used in the continuous case.
  - Matrix  $\underline{\underline{L}}$  becomes operator  $L$ .
  - Scalar product  $\underline{a}^T \underline{b}$  becomes integral  $\int a(x)b(x) dx$ .
- In somewhat loose terms,  $\nabla \chi$  is called a **sensitivity** or **Fréchet kernel** and symbolised by  $K$ .
- The only question: What is  $L^T$  in the continuous case?

## Regular wave equation

momentum balance

$$\rho(\mathbf{x})\ddot{\mathbf{u}}(\mathbf{x},t) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{x},t) = \mathbf{f}(\mathbf{x},t)$$

stress-strain relation

$$\boldsymbol{\sigma}(\mathbf{x},t) = \int_{\tau=t_0}^{\infty} \dot{\mathbf{C}}(\mathbf{x},t-\tau) : \nabla \mathbf{u}(\mathbf{x},\tau) d\tau$$

initial conditions

$$\mathbf{u}|_{t \leq t_0} = \dot{\mathbf{u}}|_{t \leq t_0} = \mathbf{0}$$

boundary conditions

$$\mathbf{n} \cdot \boldsymbol{\sigma}|_{\mathbf{x} \in \partial G} = \mathbf{0}$$

**Regular wave equation**

momentum balance

$$\rho(\mathbf{x})\ddot{\mathbf{u}}(\mathbf{x},t) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{x},t) = \mathbf{f}(\mathbf{x},t)$$

stress-strain relation

$$\boldsymbol{\sigma}(\mathbf{x},t) = \int_{\tau=t_0}^{\infty} \dot{\mathbf{C}}(\mathbf{x},t-\tau) : \nabla \mathbf{u}(\mathbf{x},\tau) d\tau$$

initial conditions

$$\mathbf{u}|_{t \leq t_0} = \dot{\mathbf{u}}|_{t \leq t_0} = \mathbf{0}$$

boundary conditions

$$\mathbf{n} \cdot \boldsymbol{\sigma}|_{\mathbf{x} \in \partial G} = \mathbf{0}$$

**Adjoint wave equation**

adjoint momentum balance

$$\rho \ddot{\mathbf{u}}^\dagger - \nabla \cdot \boldsymbol{\sigma}^\dagger = -\nabla_u \chi$$

adjoint stress-strain relation

$$\boldsymbol{\sigma}^\dagger(t) = \int_{\tau=t}^{t_1} \dot{\mathbf{C}}(\tau-t) : \nabla \mathbf{u}^\dagger(\tau) d\tau$$

terminal conditions

$$\mathbf{u}^\dagger|_{t \geq t_1} = \dot{\mathbf{u}}^\dagger|_{t \geq t_1} = \mathbf{0}$$

boundary conditions

$$\mathbf{n} \cdot \boldsymbol{\sigma}^\dagger|_{\mathbf{x} \in \partial G} = \mathbf{0}$$

notation

$$\mathbf{v} = \mathbf{u}^\dagger$$

**Regular wave equation**

momentum balance

$$\rho(\mathbf{x})\ddot{\mathbf{u}}(\mathbf{x},t) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{x},t) = \mathbf{f}(\mathbf{x},t)$$

stress-strain relation

$$\boldsymbol{\sigma}(\mathbf{x},t) = \int_{\tau=t_0}^{\infty} \dot{\mathbf{C}}(\mathbf{x},t-\tau) : \nabla \mathbf{u}(\mathbf{x},\tau) d\tau$$

initial conditions

$$\mathbf{u}|_{t \leq t_0} = \dot{\mathbf{u}}|_{t \leq t_0} = \mathbf{0}$$

boundary conditions

$$\mathbf{n} \cdot \boldsymbol{\sigma}|_{\mathbf{x} \in \partial G} = \mathbf{0}$$

**Adjoint wave equation**

adjoint momentum balance

$$\rho \ddot{\mathbf{u}}^\dagger - \nabla \cdot \boldsymbol{\sigma}^\dagger = -\nabla_u \chi$$

adjoint stress-strain relation

$$\boldsymbol{\sigma}^\dagger(t) = \int_{\tau=t}^{t_1} \dot{\mathbf{C}}(\tau-t) : \nabla \mathbf{u}^\dagger(\tau) d\tau$$

terminal conditions

$$\mathbf{u}^\dagger|_{t \geq t_1} = \dot{\mathbf{u}}^\dagger|_{t \geq t_1} = \mathbf{0}$$

boundary conditions

$$\mathbf{n} \cdot \boldsymbol{\sigma}^\dagger|_{\mathbf{x} \in \partial G} = \mathbf{0}$$

**Comments**

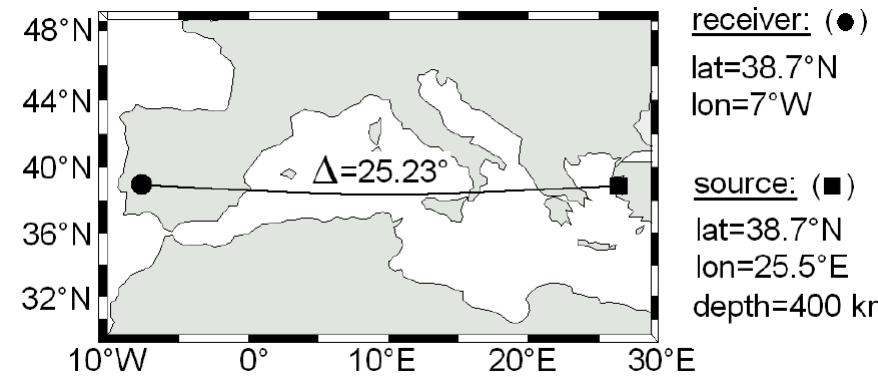
- Adjoint equation is a wave equation [same code can be used for its solution].
- Solving terminal conditions can be done by running code in reversed time.

## 4. Sensitivity kernels and kernel gallery

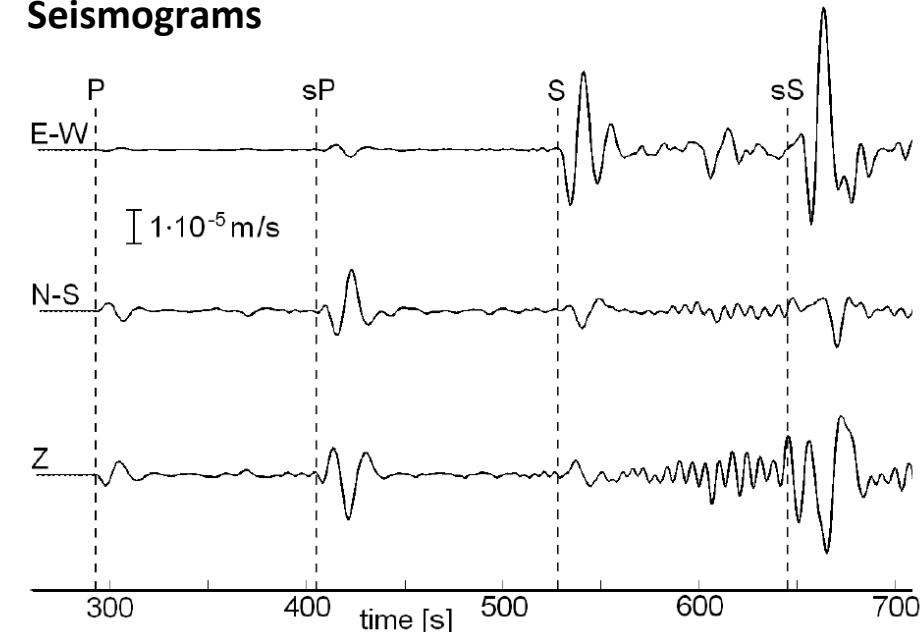
# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

## Source-receiver geometry

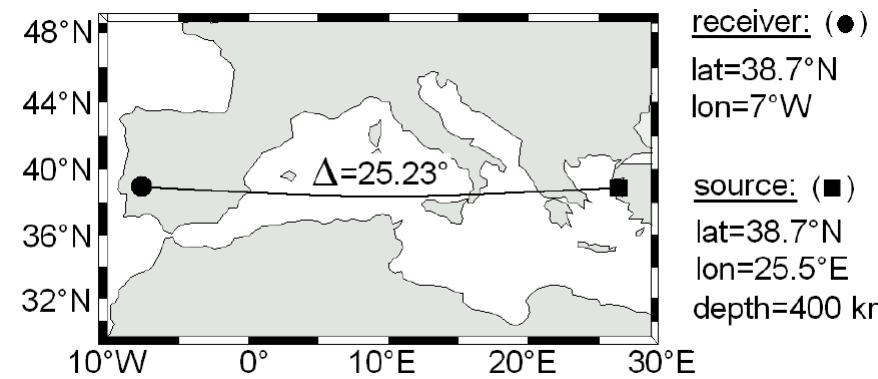


## Seismograms

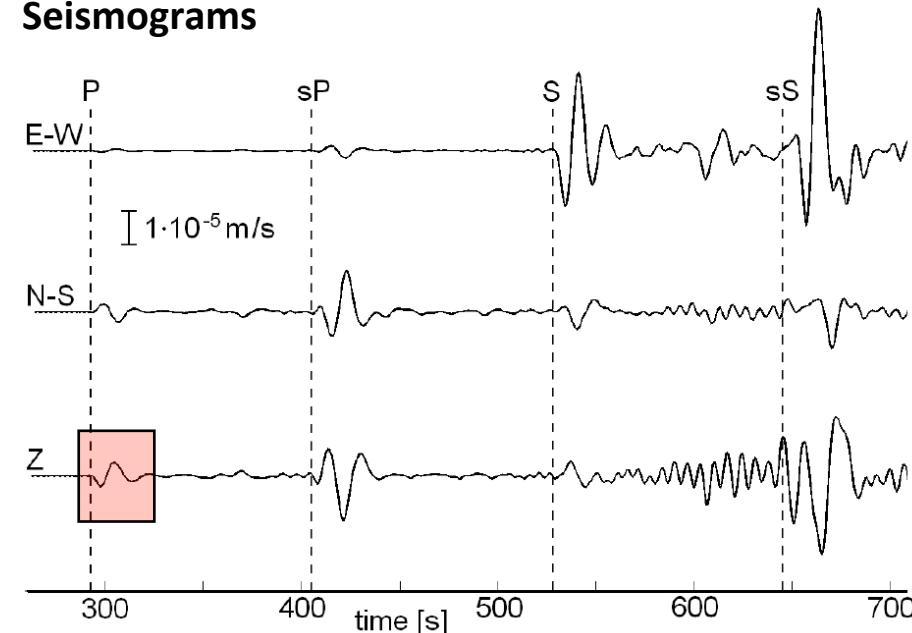


# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

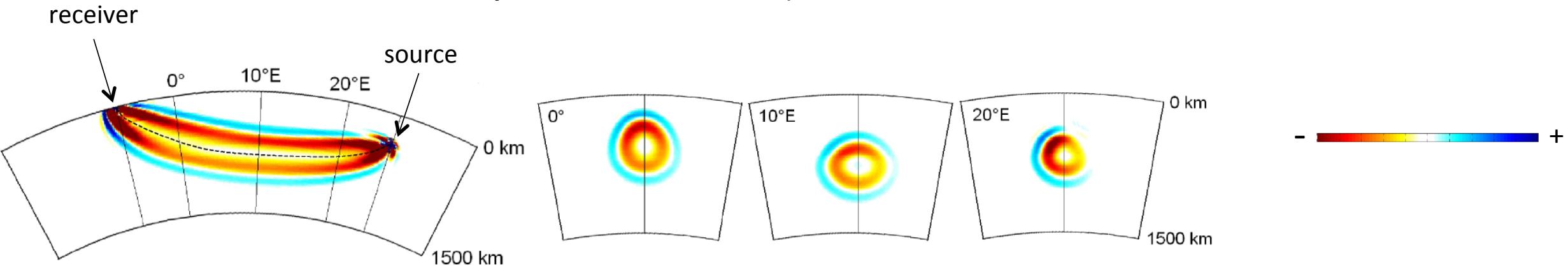
## Source-receiver geometry



## Seismograms

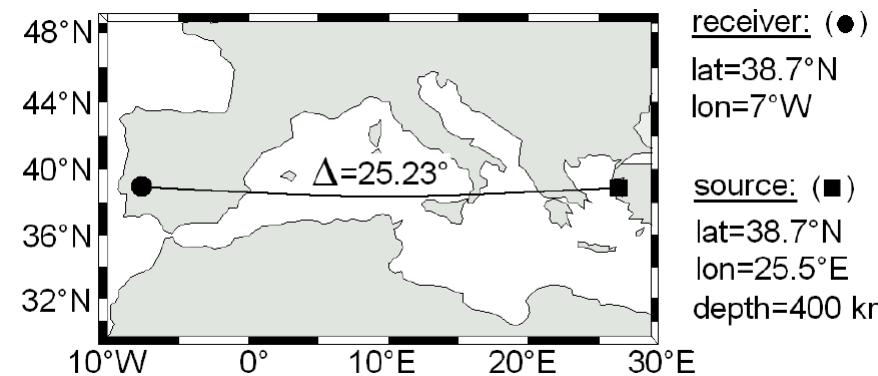


## Sensitivity kernel for P wave velocity

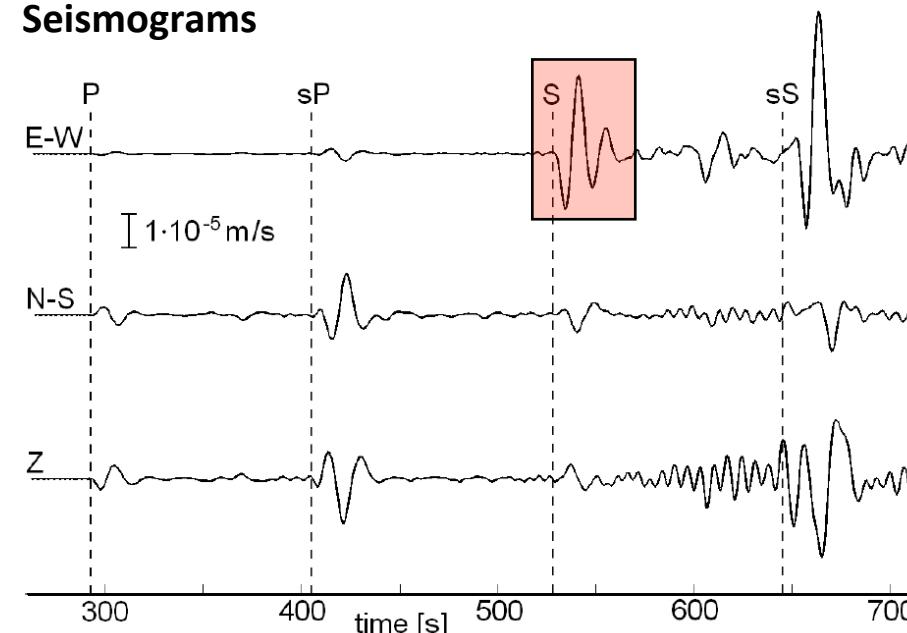


# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

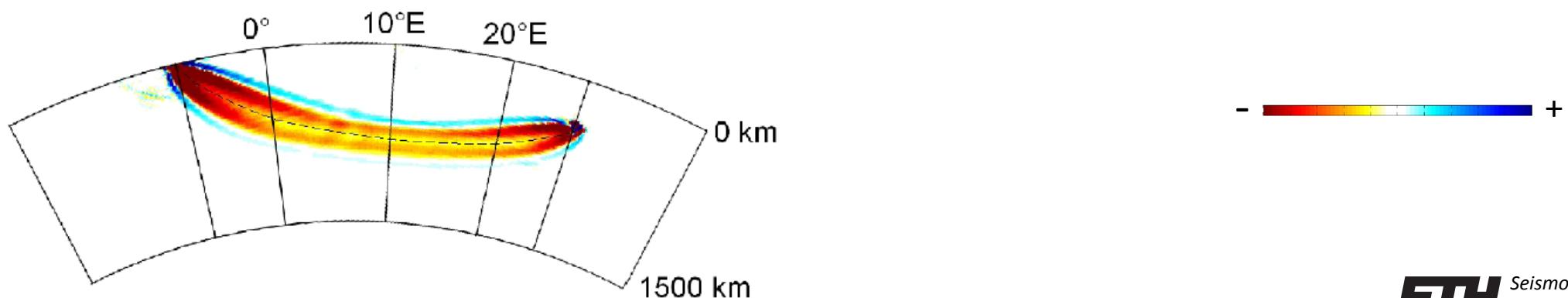
## Source-receiver geometry



## Seismograms

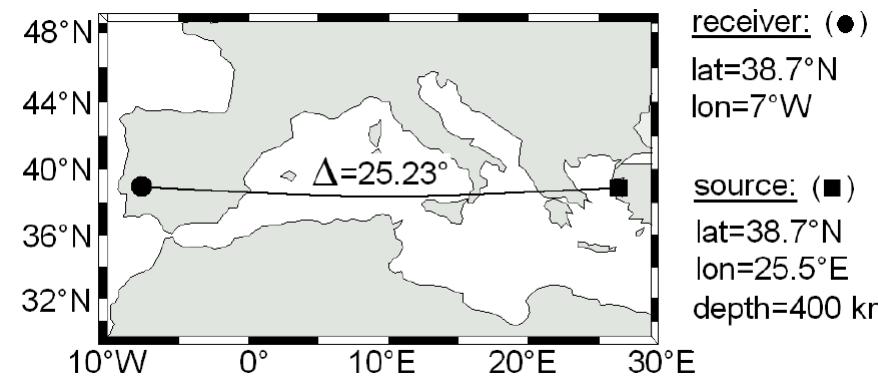


## Sensitivity kernel for S wave velocity

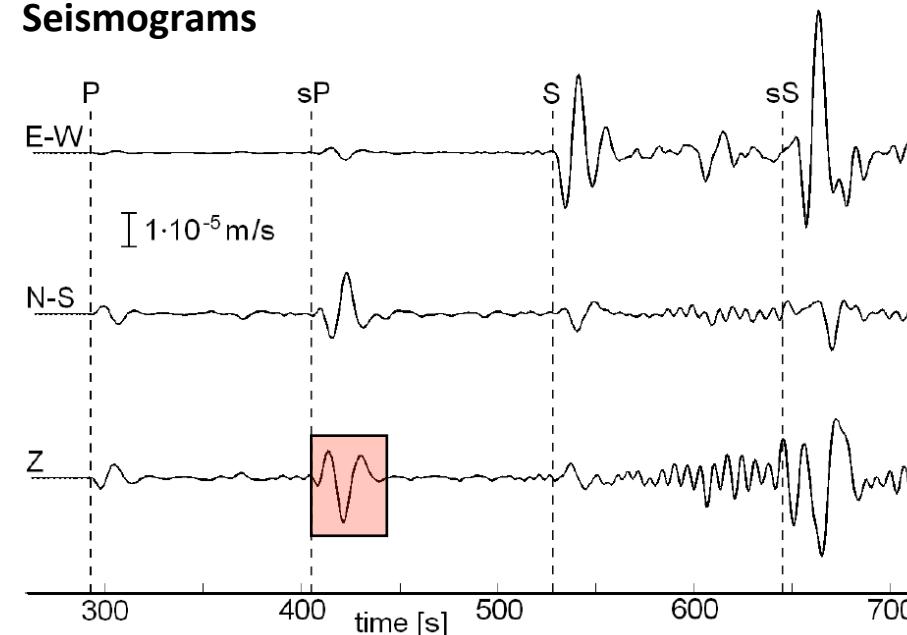


# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

## Source-receiver geometry



## Seismograms

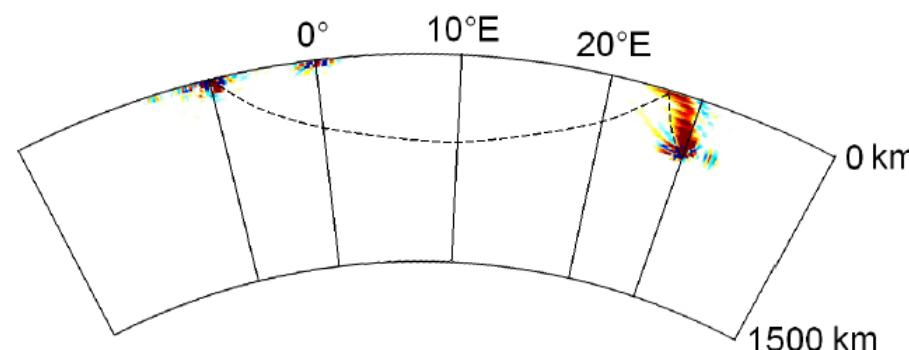
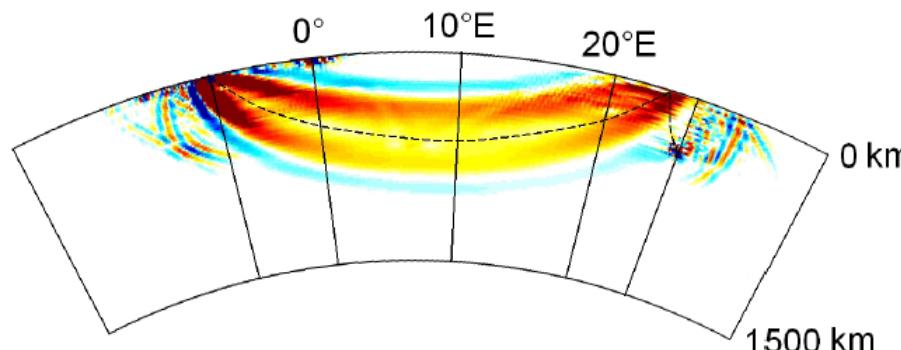


P wave velocity

Sensitivity kernels for

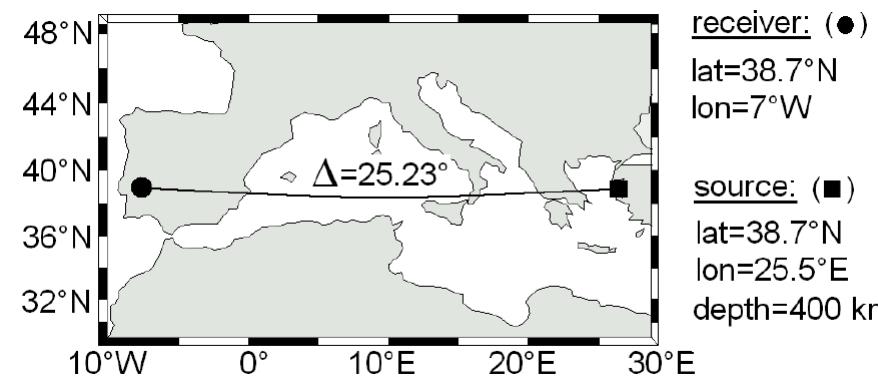
and

S wave velocity

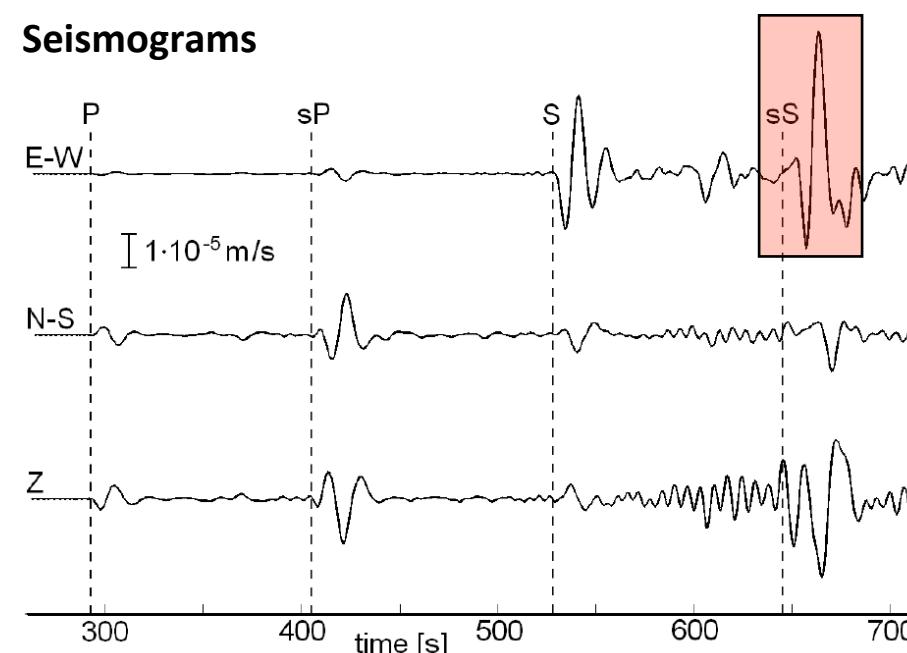


# TRAVELTIME MEASUREMENT ON SPECIFIC PHASES

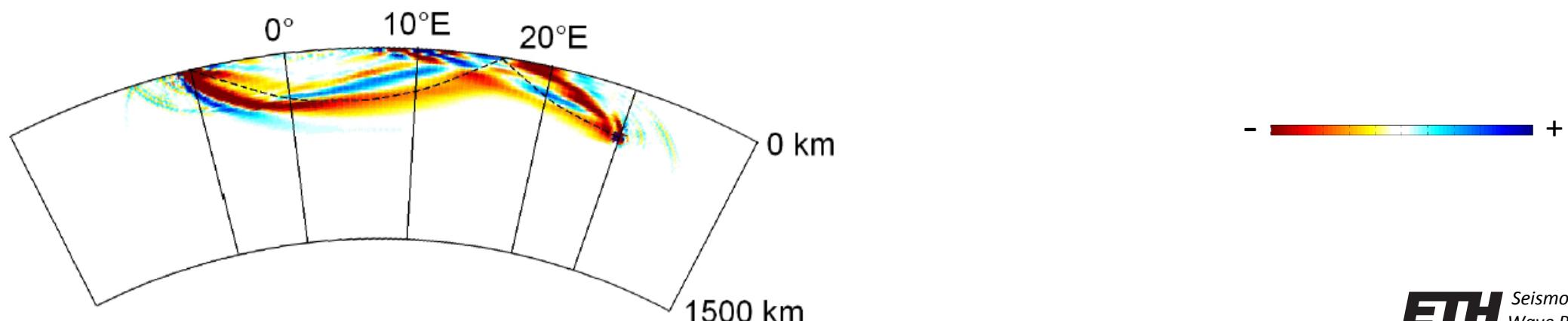
## Source-receiver geometry



## Seismograms

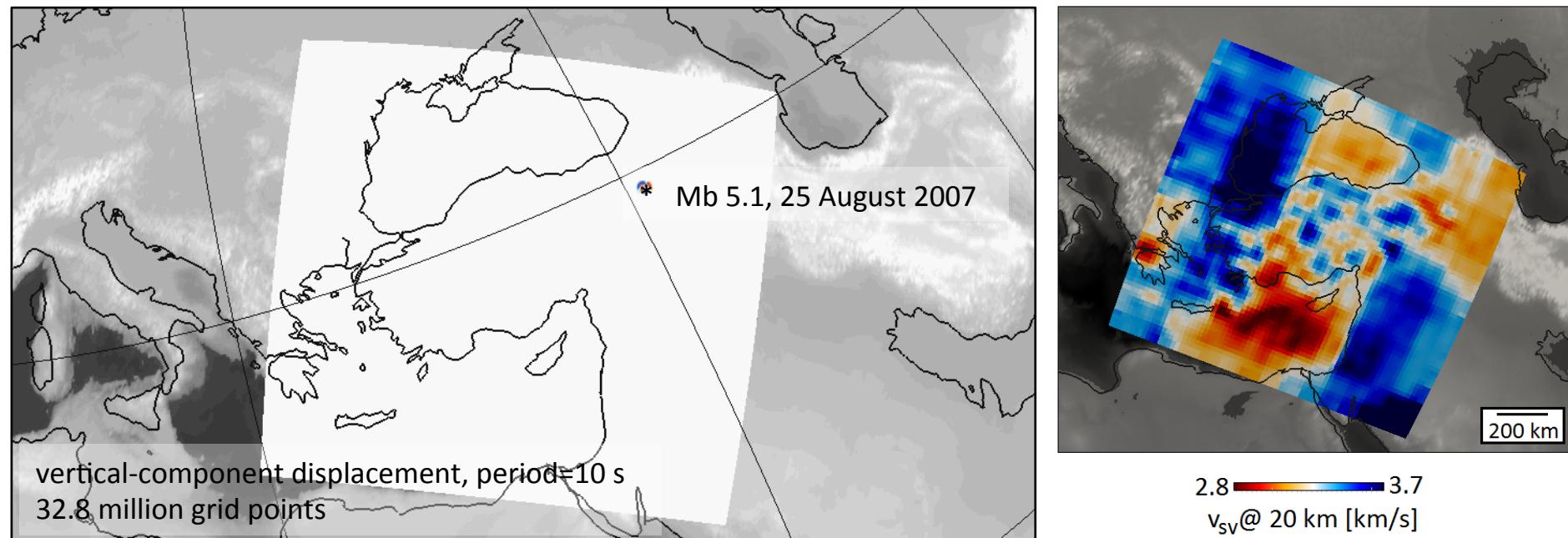


## Sensitivity kernel for S wave velocity

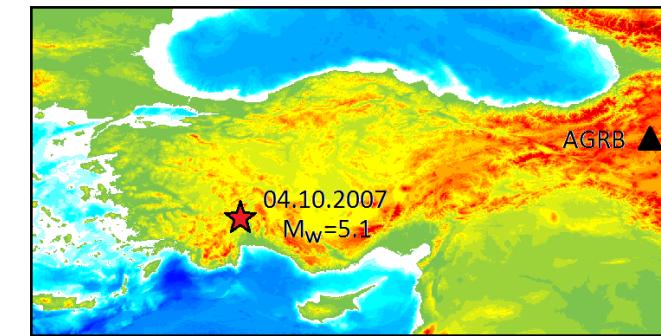
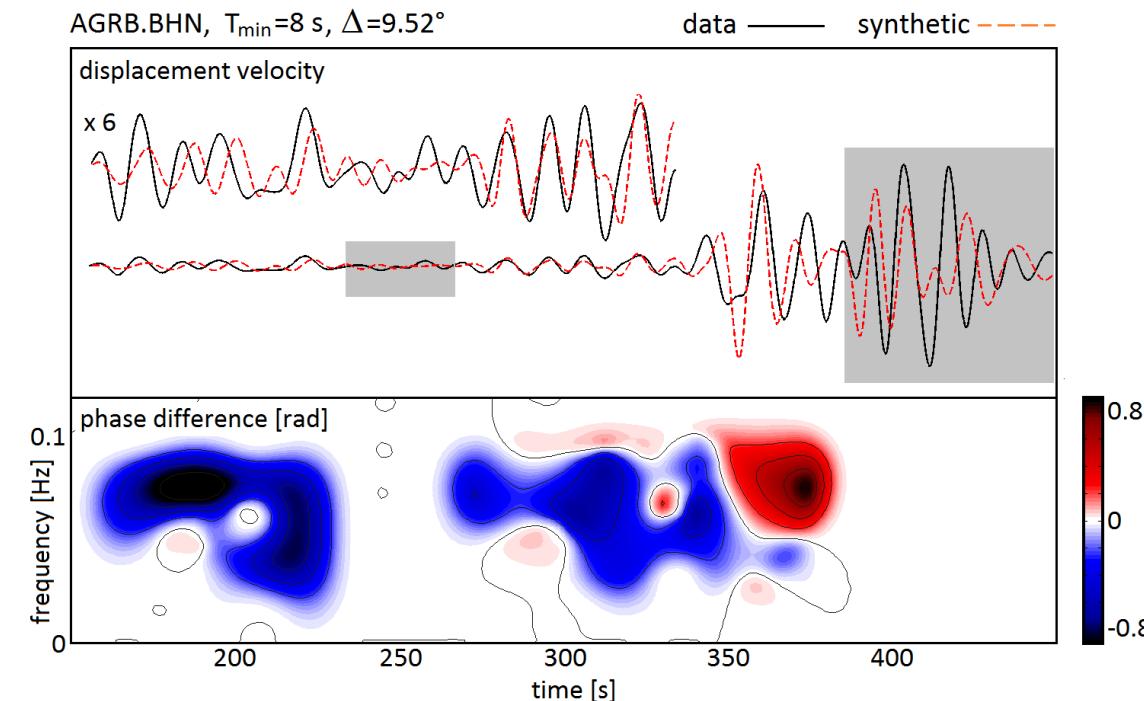


# MEASURING TIME-FREQUENCY PHASE DIFFERENCES

# MEASURING TIME-FREQUENCY PHASE DIFFERENCES

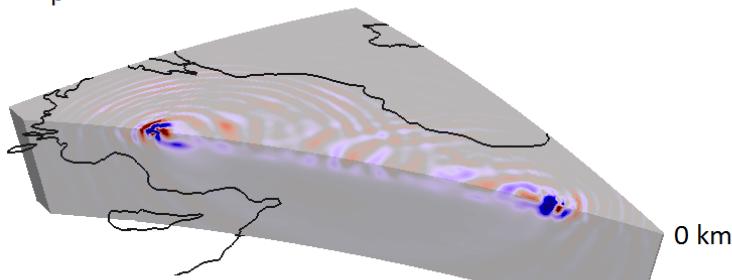
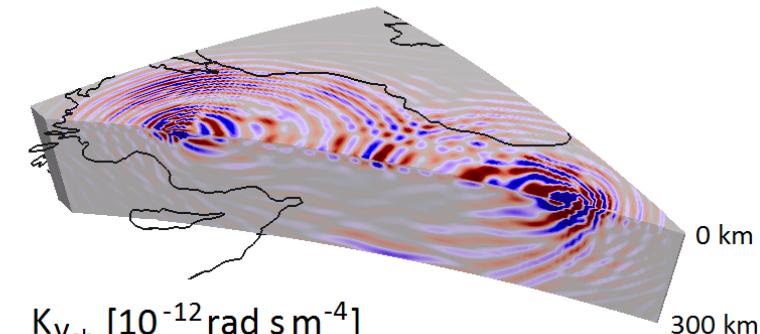
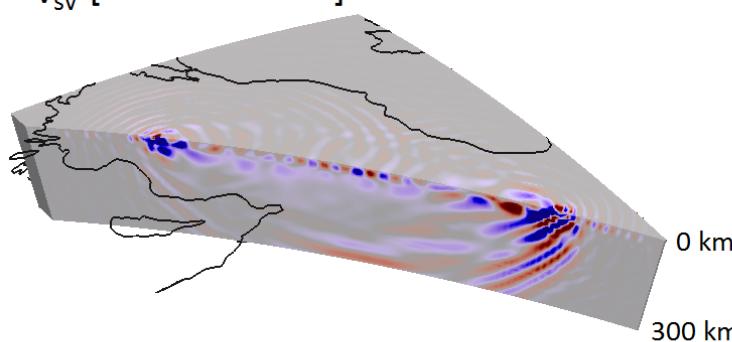
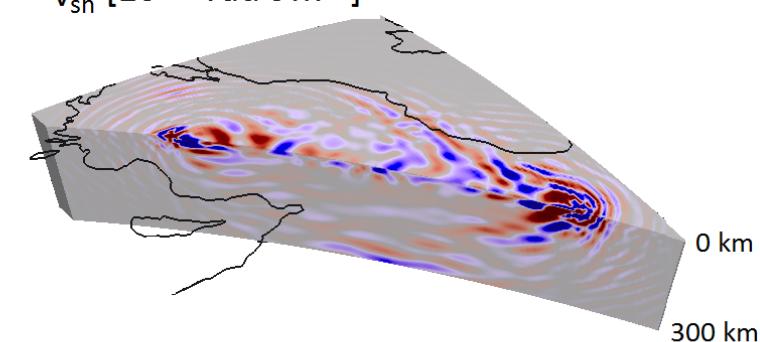


# MEASURING TIME-FREQUENCY PHASE DIFFERENCES



- Time- and frequency-dependent phase differences
- Based on selection of time windows where data and synthetics are similar
- Independent of absolute amplitudes

## Sensitivity kernels

 $K_{V_p} [10^{-12} \text{ rad s m}^{-4}]$  $K_\rho [10^{-12} \text{ rad kg}^{-1}]$  $K_{V_{sv}} [10^{-12} \text{ rad s m}^{-4}]$  $K_{V_{sh}} [10^{-12} \text{ rad s m}^{-4}]$ 

-0.4 0.4