

Lecture 1

Inverse Problems for Geophysicists: Introduction to Linear Inverse Problems and Regularization Methods

M D Sacchi
University of Alberta

Course information

- Email: msacchi@ualberta.ca
- https://github.com/msacchi/SEP_lectures

Contents

● PRELIMINARIES	3
● LINEAR INVERSE PROBLEMS	16
● TWO IMPORTANT MATRICES (OPERATORS)	36
● DISCRETE LINEAR INVERSE PROBLEMS IN EXPLICIT FORM	45
● PROBLEMS IN IMPLICIT FORM	95

PRELIMINARIES

Preliminaries: Notation

$f(\mathbf{x})$: scalar function of a vector

\mathbf{x} : a vector

x : a scalar

α : a scalar

\mathbf{M} : a matrix

$M_{i,j}$: element of the matrix

Preliminaries: ell-2 norm

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_N \end{pmatrix}$$

- Real vector

$$\mathbf{u} : N \times 1 \longrightarrow \|\mathbf{u}\|_2^2 = \mathbf{u}^T \mathbf{u} = \sum_{i=1}^N u_i^2$$

- Complex vector

$$\mathbf{u} : N \times 1 \longrightarrow \|\mathbf{u}\|_2^2 = \mathbf{u}^H \mathbf{u} = \sum_{i=1}^N u_i u_i^* = \sum_{i=1}^N |u_i|^2$$

- The ell-2 norm

$$l_2 = \|\mathbf{u}\|_2$$

Preliminaries: Matrix-times-vector

\mathbf{A} : Matrix
 \mathbf{x} : Vector
 \mathbf{y} : Vector
 α : scalar

A_{ij} : Element of the Matrix \mathbf{A}

x_i : Element of Vector \mathbf{x}

y_i : Element of Vector \mathbf{y}

Matrix Multiplication

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$
$$N \times 1 \quad N \times M \quad M \times 1$$

Preliminaries: Matrix-times-vector

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

$N \times 1$ $N \times M$ $M \times 1$

Outer Loop Inner Loop

$$y_i = \sum_{j=1}^M A_{ij} x_j \quad i = 1 \dots N$$

Output index Input index

Notation and Review of Linear Algebra

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

Matrix-times-vector multiplication ¶

```
In [1]: 1 N = 4
        2 M = 5
        3 A = randn(Float32,N,M)
        4 x = randn(Float32,M)
        5
        6 y = A*x
```

```
Out[1]: 4-element Vector{Float32}:
        -0.48016113
         3.7603946
         4.958347
         0.91909456
```

My own code...

```
In [2]: 1 y= zeros(Float32,N)
        2 for i = 1:N                                # Outer loop over data parameters
        3     for j=1:M                                # Inner loop over model parameters
        4         y[i]+=A[i,j]*x[j]
        5     end
        6 end
        7 y
```

```
Out[2]: 4-element Vector{Float32}:
        -0.48016113
         3.7603943
         4.9583473
         0.91909456
```

Preliminaries: Matrix-times-vector

$$\underset{N \times 1}{\mathbf{y}} = \underset{N \times M}{\mathbf{A}} \underset{M \times 1}{\mathbf{x}}$$

$$y_i = \sum_{j=1}^M A_{ij} x_j \quad i = 1 \dots N$$

Forward mode

$$\underset{M \times 1}{\mathbf{x}'} = \underset{M \times N}{\mathbf{A}^T} \underset{N \times 1}{\mathbf{y}}$$

$$x'_j = \sum_{i=1}^N A_{ij} y_i \quad j = 1 \dots M$$

Transpose mode

Notation and Review of Linear Algebra

$$y = A x$$

Matrix-times-vector multiplication ¶

```
In [1]: 1 N = 4
        2 M = 5
        3 A = randn(Float32,N,M)
        4 x = randn(Float32,M)
        5
        6 y = A*x
```

```
Out[1]: 4-element Vector{Float32}:
        -0.48016113
         3.7603946
         4.958347
         0.91909456
```

My own code...

```
In [2]: 1 y= zeros(Float32,N)
        2 for i = 1:N                                # Outer loop over data parameters
        3     for j=1:M                                # Inner loop over model parameters
        4         y[i]+=A[i,j]*x[j]
        5     end
        6 end
        7 y
```

```
Out[2]: 4-element Vector{Float32}:
        -0.48016113
         3.7603943
         4.9583473
         0.91909456
```

Preliminaries: Vector differentiation

- You know how to take derivatives. For instance, if $f(x) = x^2$ (**univariate calculus**) the derivative is given by

$$\frac{df(x)}{dx} = \frac{dx^2}{dx} = x \frac{dx}{dx} + \frac{dx}{dx} x = 2x$$

- What about if now we want to carry out differentiation with respect to a vector? (**Matrix/Vector derivatives or vector differentiation**). For instance

$$\frac{df(\mathbf{x})}{d\mathbf{x}} = \frac{d\|\mathbf{x}\|_2^2}{d\mathbf{x}}$$

where now \mathbf{x} is a vector.

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Preliminaries: Vector differentiation

$$\frac{d\|\mathbf{x}\|_2^2}{d\mathbf{x}} = \frac{d(\mathbf{x}^T \mathbf{x})}{d\mathbf{x}} = 2\mathbf{x}$$

Derivative of a scalar is a vector

$$\frac{d(\mathbf{u}^T \mathbf{x})}{d\mathbf{x}} = \frac{d(\mathbf{x}^T \mathbf{u})}{d\mathbf{x}} = \mathbf{u}$$

$$\frac{d(\mathbf{B} \mathbf{x})}{d\mathbf{x}} = \mathbf{B}^T$$

Derivative of a vector is a matrix

$$\frac{d(\mathbf{C}^T \mathbf{x})}{d\mathbf{x}} = \mathbf{C}$$

Preliminaries: Least-squares method via vector differentiation

$$\mathbf{y} \approx \mathbf{A}\mathbf{x}, \quad N > M \quad \mathbf{A} \in \mathcal{R}^{N \times M}$$

$$\epsilon = \|\mathbf{e}\|_2^2 = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad \text{Error function to minimize}$$

$$\begin{aligned} \epsilon &= \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{A}\mathbf{x})^T (\mathbf{y} - \mathbf{A}\mathbf{x}) \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{A}\mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} \end{aligned}$$

$$\begin{aligned} \frac{d\epsilon}{d\mathbf{x}} &= \frac{d}{d\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} + \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{A}\mathbf{x}] \\ &= 2\mathbf{A}^T \mathbf{A}\mathbf{x} - 2\mathbf{A}^T \mathbf{y} = 0 \end{aligned}$$

Preliminaries: Least-squares method via vector differentiation

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y} \quad \xrightarrow{\text{Assuming invertibility}}$$

$$\mathbf{x}_{sol} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

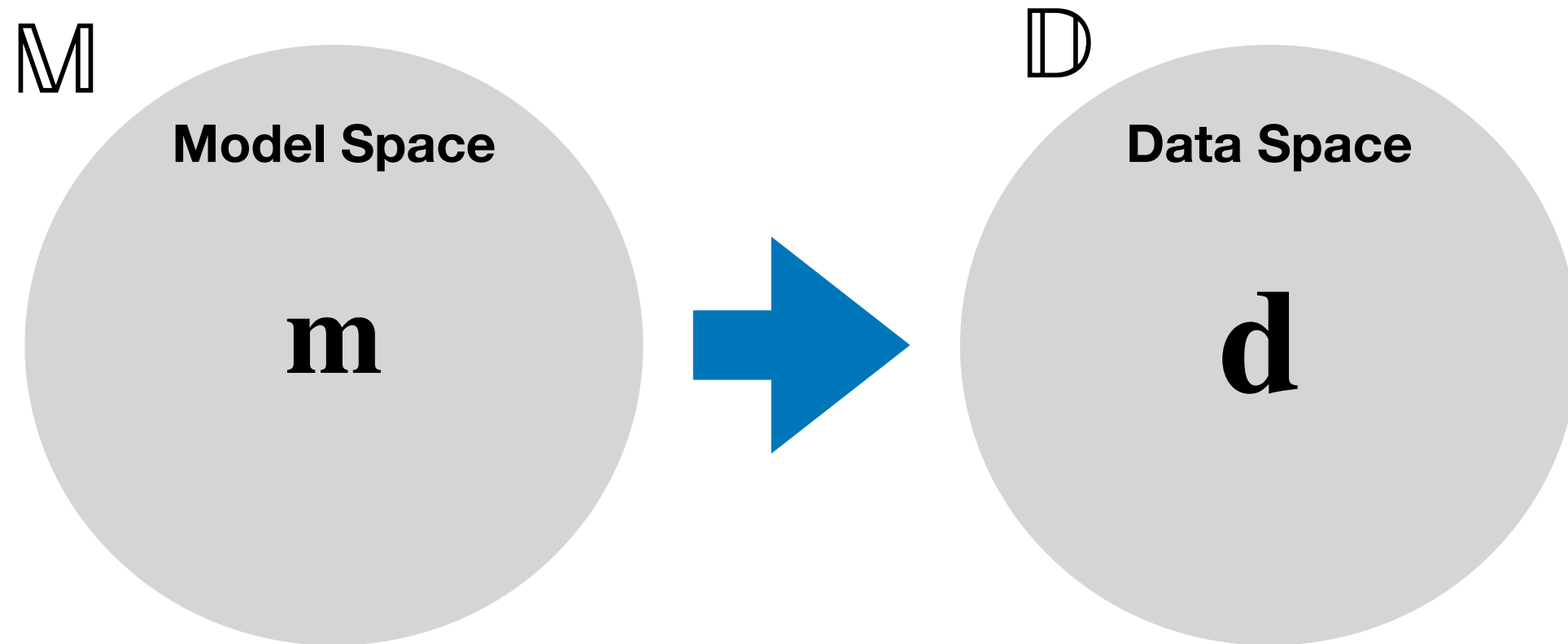
LINEAR INVERSE PROBLEMS

Inverse Problems

- An Inverse Problem is a mathematical problem where one attempts to estimate **models** that explain **observations**. We often name the observations **d (data)**
- Observations are generally measured on the surface of the earth and are discrete in time and space
- The subsurface is described by **properties** (density, velocity, reflectivity, resistivity, etc). These properties exist everywhere. We will refer to these properties as **$m(x,y,z)$** or **\mathbf{m} (model parameters)**
- **$m(x,y,z)$** is the distribution of some property that can be discretized and represented via the vector **\mathbf{m}**

Two spaces

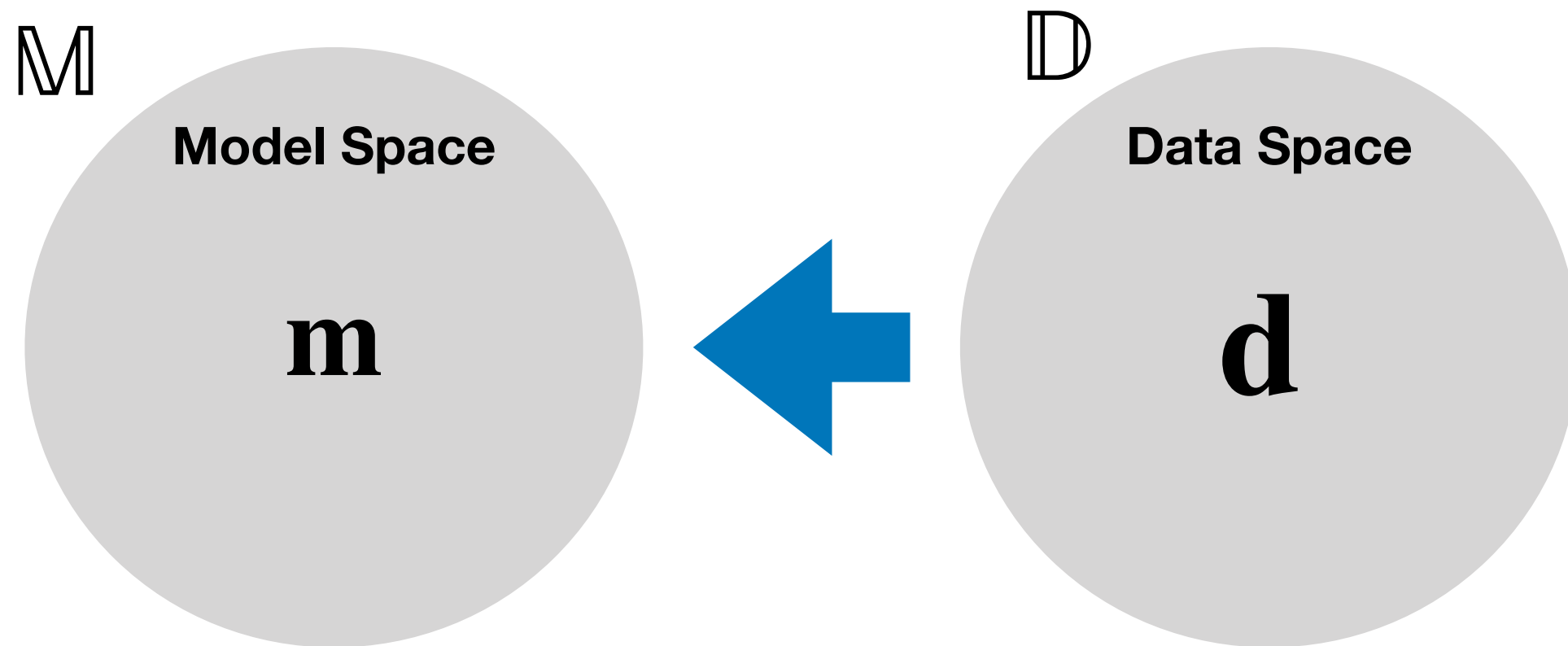
- Forward Problem



$$F[\mathbf{m}] = \mathbf{d}$$

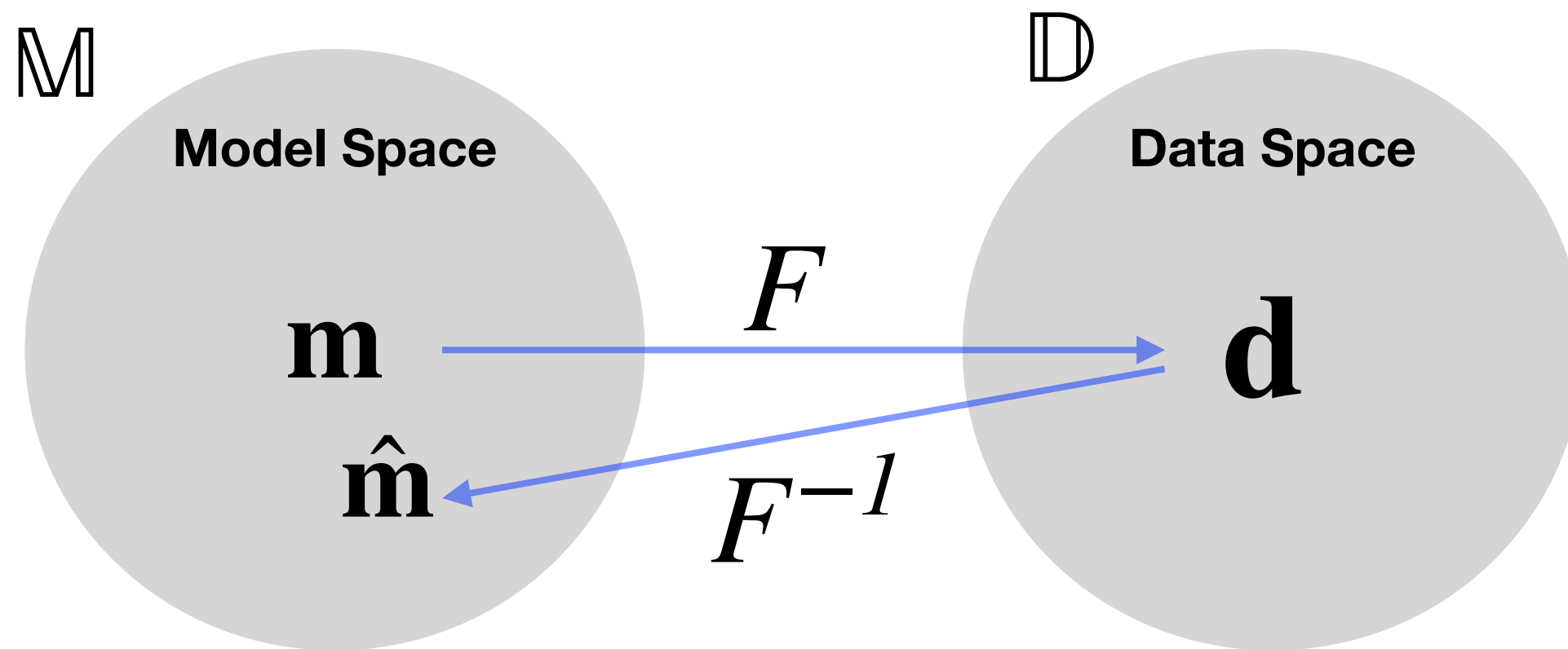
Two spaces

- Inverse Problem



$$F^{-1}[\mathbf{d}] = \mathbf{m}$$

It is more complicated...



$\hat{\mathbf{m}}$: Solution

Geophysics and Inverse Problems

Data (What you can measure)	Model (What you would like to know)	Method
Gravity anomalies	Density	Potential Field Methods
Electrical potential	Resistivity	Potential Field Methods
Electrical and Magnetic Field	Electrical conductivity	EM/MT methods
Magnetic Fields	Susceptibility	Potential Field Methods
Seismic Waves	Velocities	Seismic Methods

ILL-posed problems

- **Well-posed problem.** A problem is said to be well-posed when
 - There is a solution
 - The solution is unique
 - The solution is stable
- If one of the above is not true, the problem is called an **ill-posed problem**
- Typical geophysical inverse problems are ill-posed problems (2 and 3 are not true)



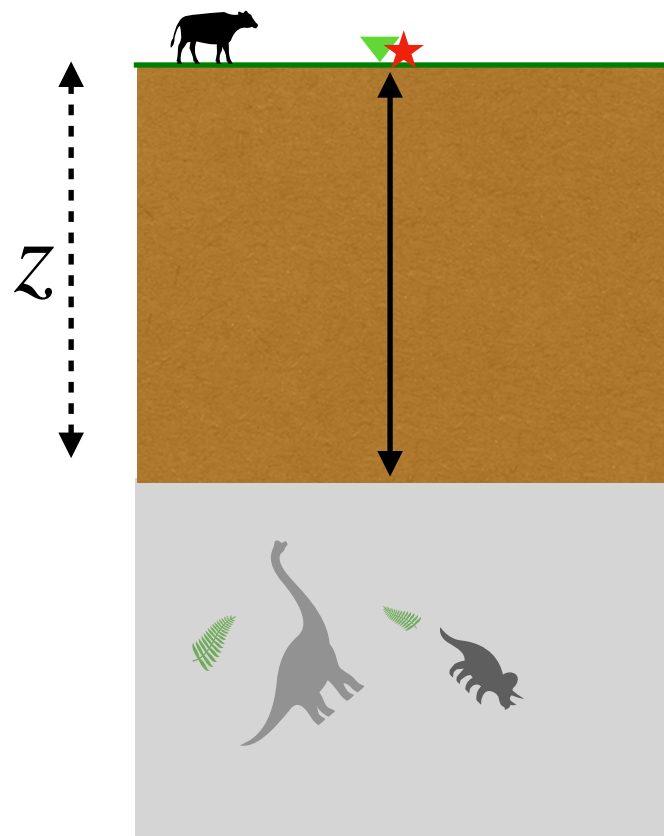
Jacques Hadamard, 1865-1963

ILL-posed problems

- There is a solution: YES
 - There is a solution otherwise we wouldn't be here
 - e.g. Rocks have density causing gravity anomalies

ILL-posed problems

- The solution is unique: **Generally NO**
- For instance, Depth-Velocity ambiguity: $t = \frac{2z}{v}$
- Assume two-way traveltime for a vertical incidence plane wave



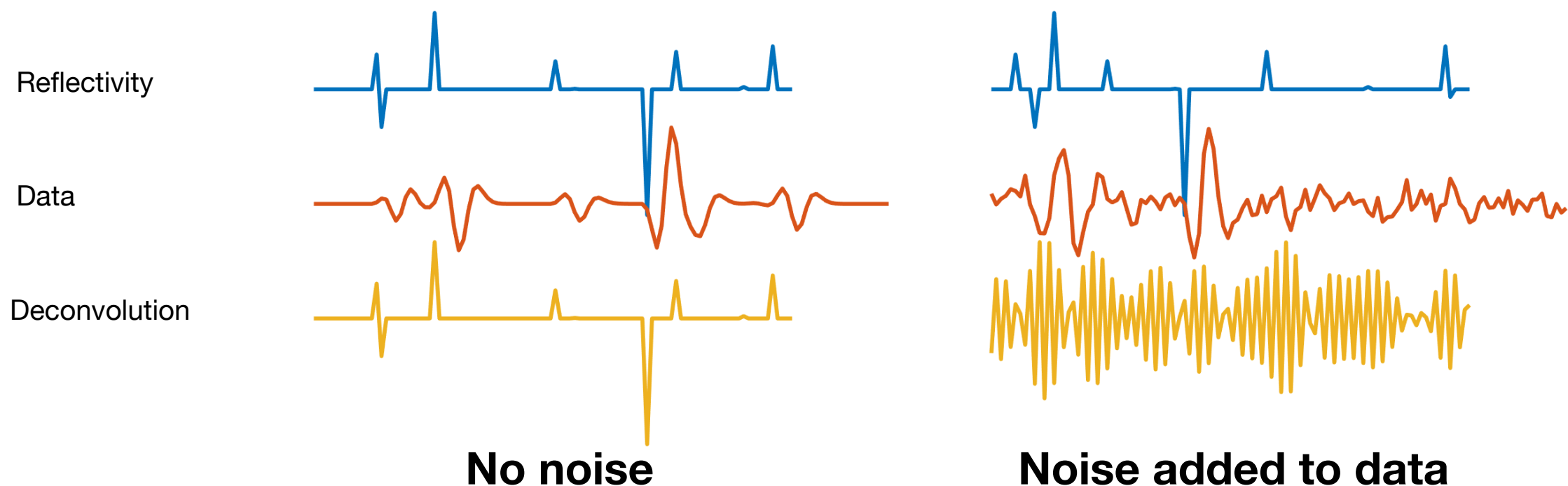
$$t = \frac{2 \times 1000m}{2000m/s} = \frac{2 \times 1500m}{3000m/s}$$

Sol 1	Sol 2
1000m	1500m
2000m/s	3000m/s

Clearly, there is an infinite number of solutions

ILL-posed problems

- The solution is stable: **Not true**
- A small perturbation in the data causes a large perturbation in the solution
- e.g. Deconvolution



ILL-posed problems

- The solution is stable: **Not true**
- A small perturbation can cause a large perturbation in the solution
- Assume **perfect** data and an invertible operator

$$d = F(m) \rightarrow \hat{m} = F^{-1}(d)$$

- Assume perturbed data (inaccurate data)

$$\hat{m} + \delta m = F^{-1}(d + n)$$

ILL-posed problems

- The solution is stable: **Not true**
- Solution for perturbed data

$$\hat{m} + \delta m = F^{-1}(d + n)$$

- The solution is not stable. The latter means

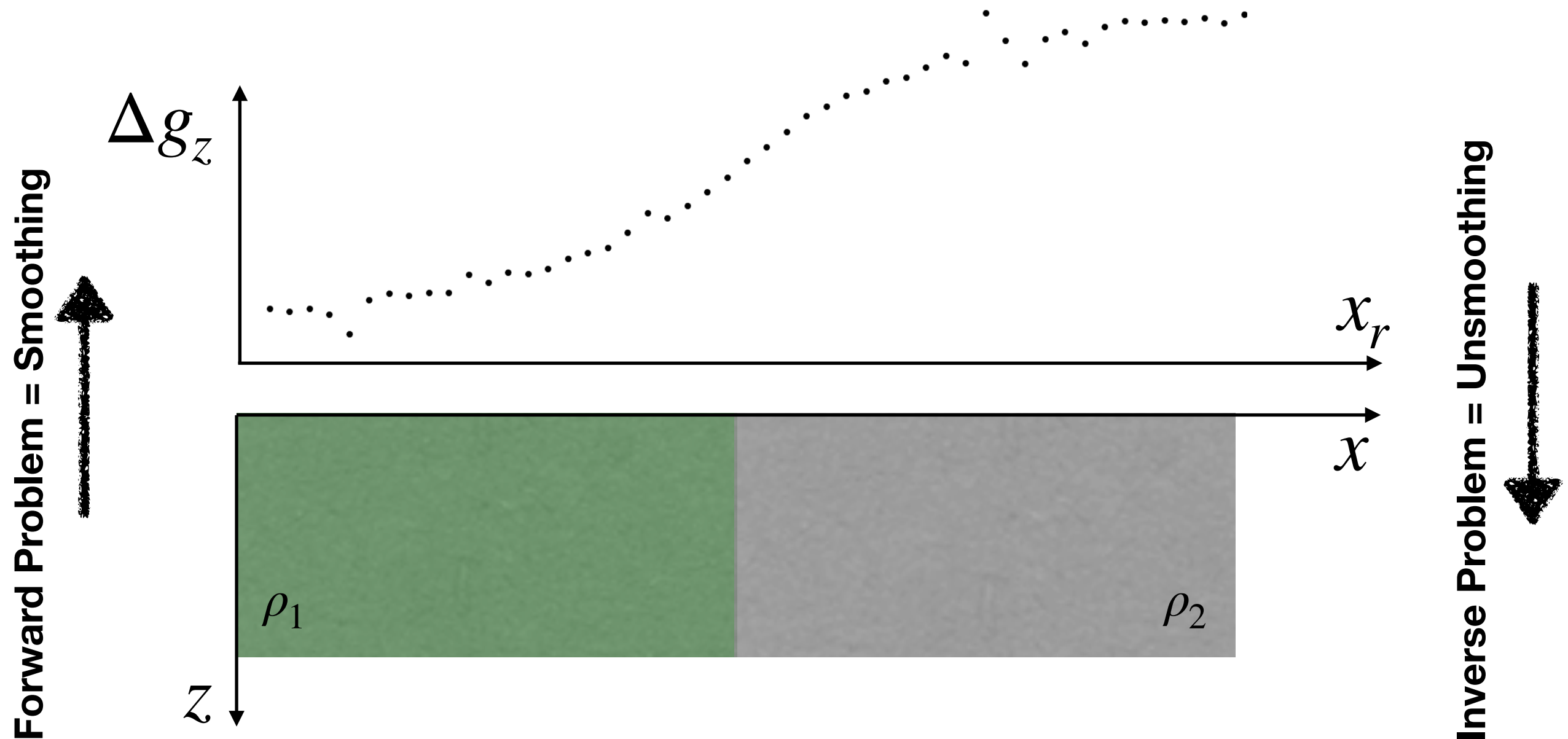
$$\text{if } |n| = \text{small}, |\delta m| = \text{large}$$

- This is the main problem we face in geophysical inversion and often in data processing. You can think that somewhere one is trying to divide by small numbers and hence, whatever you are doing to your data leads to “noise amplification”.

Why geophysical problems are ill-posed?

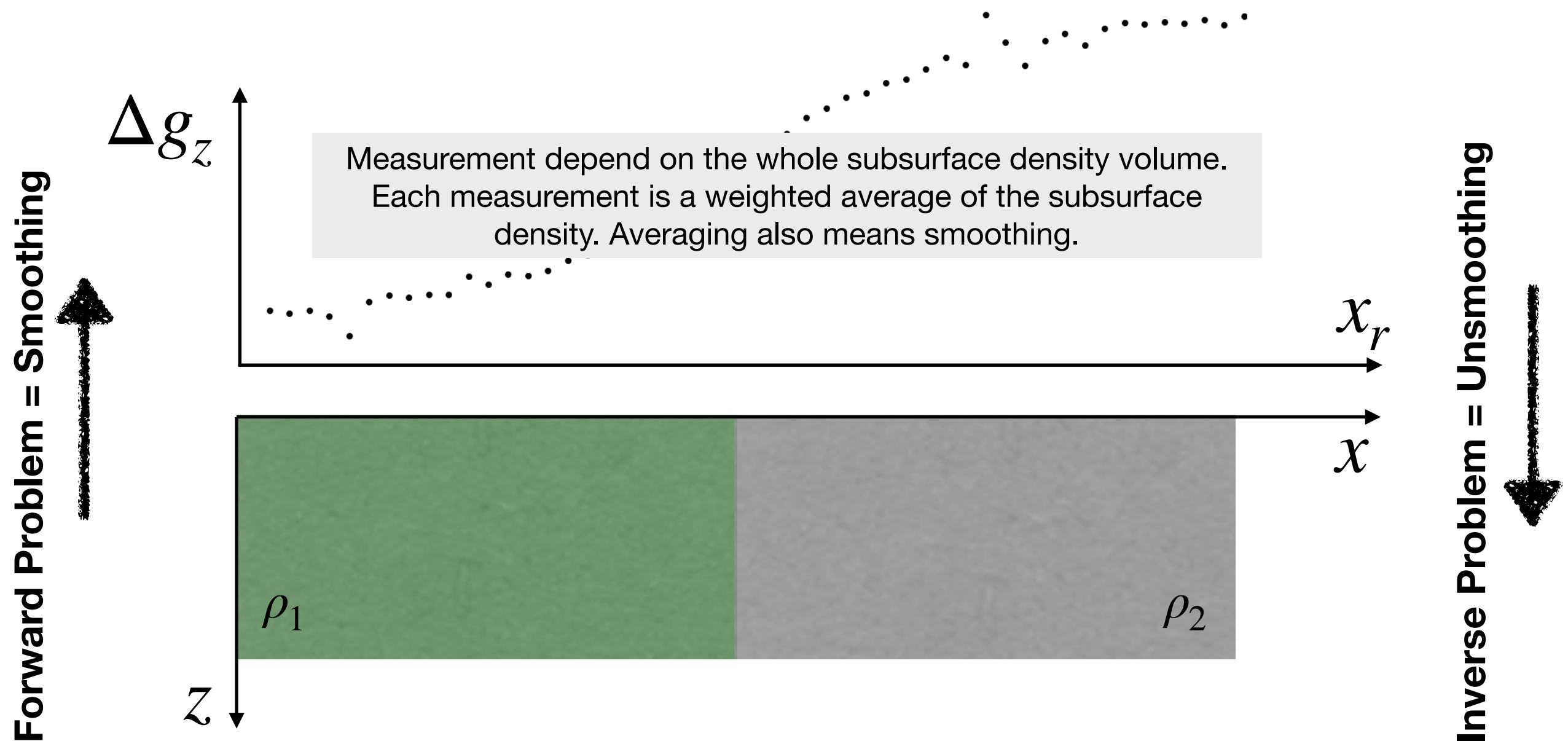
- Not enough data (insufficient spatial data or insufficient bandwidth)
- Presence of noise also leads to unstable solutions
- A central characteristic of the forward problem:
 - The Forward Problem smooths the unknown subsurface properties. Seismic waves, gravity anomalies, electrical potential, etc are nice smooth functions of space and time. Subsurface properties, on the other hand, might or might not be smooth!
- The next slide illustrates the aforementioned concept

Example: Inversion of gravity anomalies



$$\Delta g_z(x_r, z_r = 0) = \int_{x,y} G(x, y | x_r, z_r = 0) \Delta \rho(x, y) dx dy$$

Example: Inversion of gravity anomalies



$$\Delta g_z(x_r, z_r = 0) = \int_{x,y} G(x, y | x_r, z_r = 0) \Delta \rho(x, y) dx dy$$

Key points

- Previous example shows an interesting feature of the Forward Problem:
 - Observations (data) are a weighted average of density. The averaging kernel smooths the density anomaly (The physics, in this case given by Newton's gravitational attraction, is the cause of this behaviour)
 - Inverse Problem: Recovering the density from the data entails the opposite of smoothing (un-smoothing); an unstable operation
 - *Smoothing* = is a low pass operation = **Stable**
 - *Un-smoothing* = is a high pass operation = **Unstable**
- Solving an inverse problem in many cases entails controlling instability

Key points

- Subsurface properties (such as density anomalies) are a function of space (that exists everywhere). The observed gravity anomalies are given by a finite number of observations

$$d(\mathbf{r}_j) = \int G(\mathbf{x} | \mathbf{r}_j) m(\mathbf{x}) d\mathbf{x} \quad j = 1 \dots N$$

- So inversion, theoretically, attempt to go from a finite number of observations to a function

$$\mathbf{d} \in R^N, m(\mathbf{x}) \in R^\infty$$

- There is a natural **underdeterminacy** in an inverse problem

Linear inverse problems

- Many inverse problems in seismology entail solving integral equations. Others entail solving PDEs. We will start with simple problems that can be written as integral equations that after discretization lead to discrete system of equations.
 - Integral equations
 - Linear discrete system of equations
 - Regularization methods
 - Connection to exploration seismology:
 - Deconvolution
 - AVO inversion
 - Linearized seismic imaging (Forward and Adjoint Operators)
 - Migration & Least-Squares migration

Linear inverse problems

- Fredholm integral equation of 1st kind

$$d(r_j) = \int_X G(r_j, x) m(x) dx \quad j = 1 \dots N$$

- We can discretize the model

$$x_k = x_0 + (k - 1)\Delta x \quad k = 1 \dots M$$

$$d(r_j) = \sum_{k=1}^M G(r_j, x_k) m(x_k) \Delta x \quad j = 1 \dots N$$

- Which leads to a system of equations

$$\mathbf{d} = \mathbf{Lm}, \quad \mathbf{d} \in R^N, \quad \mathbf{m} \in R^M$$

Examples

$$d(r_j) = \int_X G(r_j, x) m(x) dx \quad j = 1 \dots N$$

$$s(t_i) = \int w(t_i - \tau) r(\tau) d\tau$$

Convolution

$$T(r_j) = \int_0^{r_j} s(l) dl$$

Travel-time tomography

$$d(\omega, s_i, r_j) = \int_x \int_z B(\omega, x, y | s_i, r_j) m(x, z) dx dz$$

Born imaging

TWO IMPORTANT MATRICES (OPERATORS)

Forward and Adjoint Operators

- Two special operations (or operators?)

Forward : $\mathbf{d} = \mathbf{L}\mathbf{m}$

Transpose or adjoint : $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{d}$

- \mathbf{L} is a linear operator (Forward modelling operator).
- Why are them special?
 - Iterative solvers for large inverse problems only need to evaluate $\mathbf{L}[\cdot]$ and $\mathbf{L}^T[\cdot]$
We pretend these are Matrices but in reality these are linear operators (codes)
 - I usually interpret operators as matrices. In reality, operators are codes that apply an operation *on the flight* (implicit form)
 - *We will see these operators everywhere today*

Forward and Adjoint Operators

- **When \mathbf{L} is a real matrix**

$$\mathbf{L}' = \mathbf{L}^T$$

- When \mathbf{L} is a complex matrix

$$\mathbf{L}' = \mathbf{L}^H$$

Forward and Adjoint Operator

- Two special operations (or operators)

$$\text{Forward : } \mathbf{d} = \mathbf{L}\mathbf{m} \quad (1)$$

$$\text{Transpose or adjoint : } \tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{d} \quad (2)$$

- Replace (1) into (2)

$$\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{L}\mathbf{m}$$

- Questions:
 - Is the *adjoint model* a good representation of the *true model* ?
 - Can I remove the distortion ?

See: Migration vs. Least-squares Migration: *Least-squares migration of incomplete reflection data. Nemeth et al. GEOPHYSICS (1999), 64(1)*

Forward and Adjoint Operators

- **Two special operations (or operators)**

Forward : $\mathbf{d} = \mathbf{L}\mathbf{m}$ (1)

Conjugate transpose or adjoint : $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{d}$ (2)

(2) into (1) $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{L}\mathbf{m}$

$\tilde{\mathbf{m}}$ is a distorted version of \mathbf{m}

Forward and Adjoint Operators

Forward : $\mathbf{d} = \mathbf{L}\mathbf{m}$

Conjugate transpose or adjoint : $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{d}$

$$\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{L} \mathbf{m}$$

$$\tilde{\mathbf{m}} \approx \mathbf{m} \text{ if } \mathbf{L}^T \mathbf{L} \approx \mathbf{I}$$

The adjoint, in some instances, is good enough to estimate \mathbf{m} . An example is RTM which can be considered the adjoint of Born modelling

Forward and Adjoint Operators

Forward : $\mathbf{d} = \mathbf{L}\mathbf{m}$

Conjugate transpose or adjoint : $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{d}$

$$\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{L} \mathbf{m}$$

$$\mathbf{L}^T \mathbf{L} \approx \mathbf{diag}(\mathbf{L}^T \mathbf{L}) = \mathbf{D}, \quad \hat{\mathbf{m}} = \mathbf{D}^{-1} \tilde{\mathbf{m}}$$

This is an inexpensive way of partially correcting the adjoint to make it look closer to the desired solution

Forward and Adjoint Operators

Example: **Averaging**

$$d = [1/3 \ 1/3 \ 1/3] \begin{bmatrix} m_1 \\ m_2 \\ m_3 \end{bmatrix} = \mathbf{L}\mathbf{m} \quad \text{Averaging operator}$$

$$\tilde{\mathbf{m}} = \begin{bmatrix} \tilde{m}_1 \\ \tilde{m}_2 \\ \tilde{m}_3 \end{bmatrix} = \begin{bmatrix} d/3 \\ d/3 \\ d/3 \end{bmatrix} = \mathbf{L}^T d \quad \text{Adjoint of the averaging operator}$$

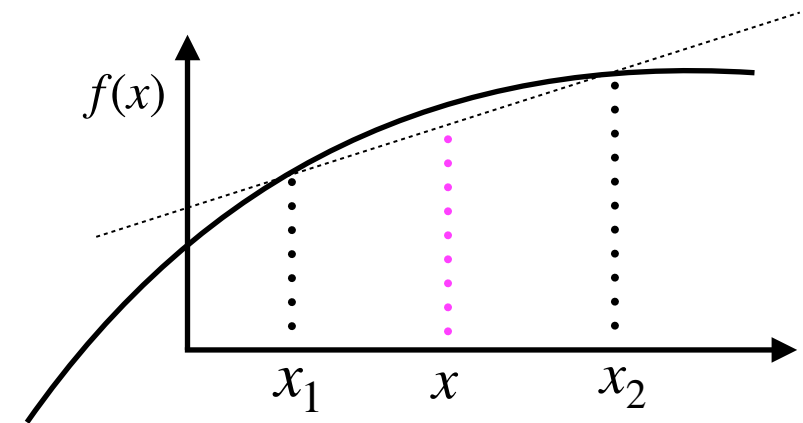
Question: Can you recover \mathbf{m} from $\tilde{\mathbf{m}} = \mathbf{L}^T \mathbf{L}\mathbf{m}$?

Forward and Adjoint Operators

Example: **Linear Interpolator (this is used a lot!!)**

$$f(x) = a f(x_1) + b f(x_2)$$

Linear Interpolator (LI)



$$f(x) = \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} f(x_1) \\ f(x_2) \end{bmatrix}$$

$$\begin{bmatrix} \tilde{f}(x_1) \\ \tilde{f}(x_2) \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} f(x)$$

$$\begin{aligned} \tilde{f}(x_1) &= a f(x) \\ \tilde{f}(x_2) &= b f(x) \end{aligned}$$

Adjoint of LI

Forward: 2 samples make 1, Adjoint: 1 sample makes 2

DISCRETE LINEAR INVERSE PROBLEMS IN EXPLICIT FORM (WITH MATRICES)

Discrete Linear Inverse Problems

$$\mathbf{d} = \mathbf{L}\mathbf{m}, \quad \mathbf{L} \in \mathcal{R}^{N \times M}$$

- A) $N > M$: More observations than data = Overdetermined problem
- B) $N < M$: Less observations than data = Underdetermined problem
- Let's consider Problem B)
 - Underdetermined Problem with Accurate Data
 - Underdetermined Problem with Noisy (inaccurate) Data

Discrete Linear Inverse Problems

(1) Accurate data $\mathbf{d} = \mathbf{Lm}$

(2) Inaccurate data $\mathbf{d} = \mathbf{Lm} + \mathbf{e} \rightarrow \mathbf{d} \approx \mathbf{Lm}$

**Consider discrete data and discrete
vector of model parameters**

Discrete Linear Inverse Problem with Accurate Data

$$\mathbf{d} = \mathbf{L}\mathbf{m}$$

$$\mathbf{d} \in R^N, \quad \mathbf{m} \in R^M$$

We assume an under-determined problem $M > N$

- We have more unknowns than observations.
- Hence, we have an infinite number of solutions.
- One needs to pick one. Which one?

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- Among all possible solution select the one with minimum Euclidian norm. In other words,

$$\textbf{Minimize } \|\mathbf{m}\|_2^2$$

$$\textbf{Subject to } \mathbf{d} = \mathbf{Lm}$$

- The latter is solved using the Method of Lagrange Multipliers. We minimized an augmented cost function

$$J = \|\mathbf{m}\|_2^2 + \mathbf{b}^T(\mathbf{Lm} - \mathbf{d})$$

- $\mathbf{b} \in R^N$ is the vector of Lagrange Multipliers, one multiplier per observation

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- To solve the problem we need to minimize the cost respect to \mathbf{m} and \mathbf{b}

$$J = \|\mathbf{m}\|_2^2 + \mathbf{b}^T(\mathbf{L}\mathbf{m} - \mathbf{d})$$

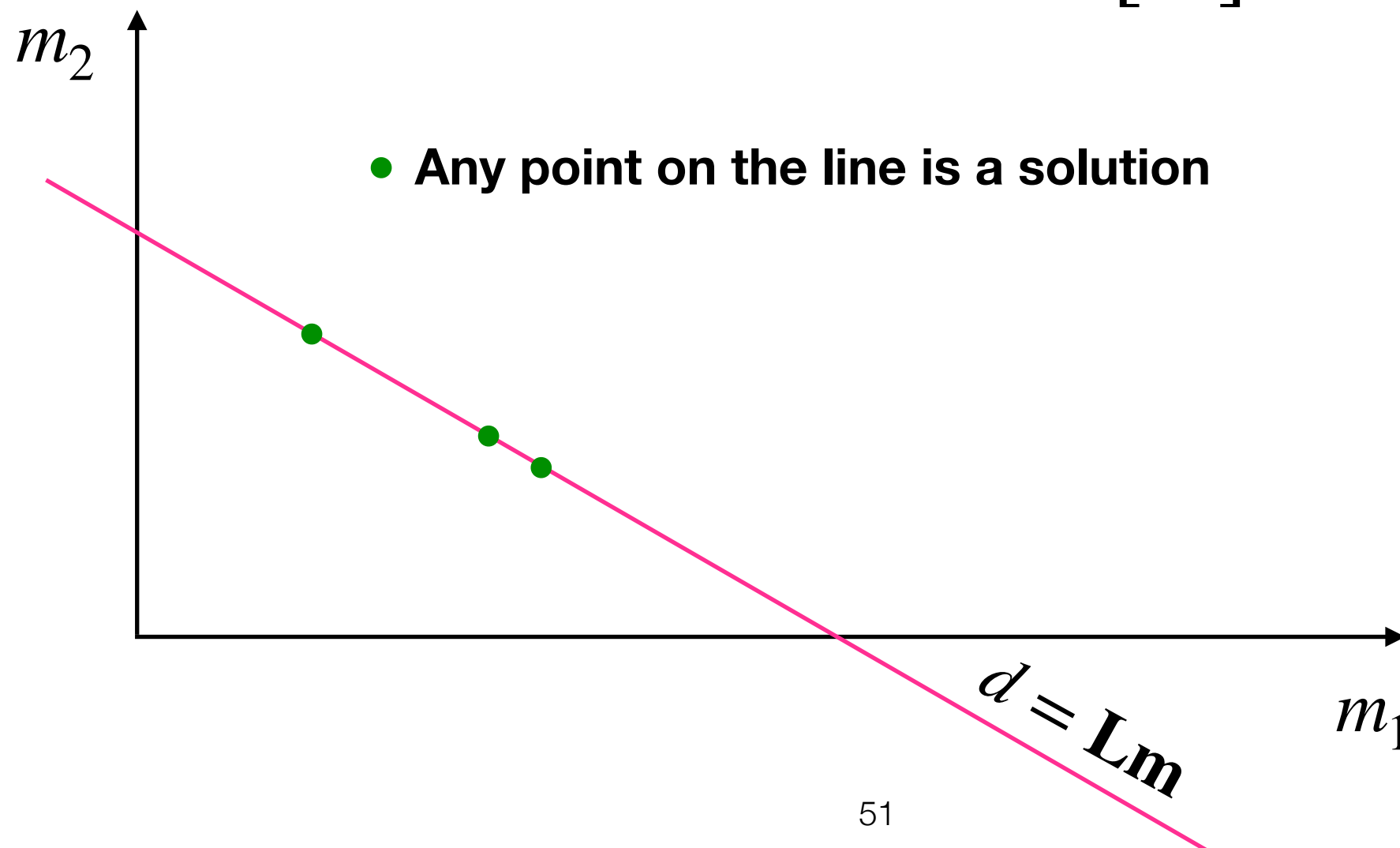
- Which leads to the so called minimum norm solution (I will do it in the whiteboard)

$$\mathbf{m}_{mn} = \mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{d}$$

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- $N=1, M=2$ (2 unknowns and 1 equation)

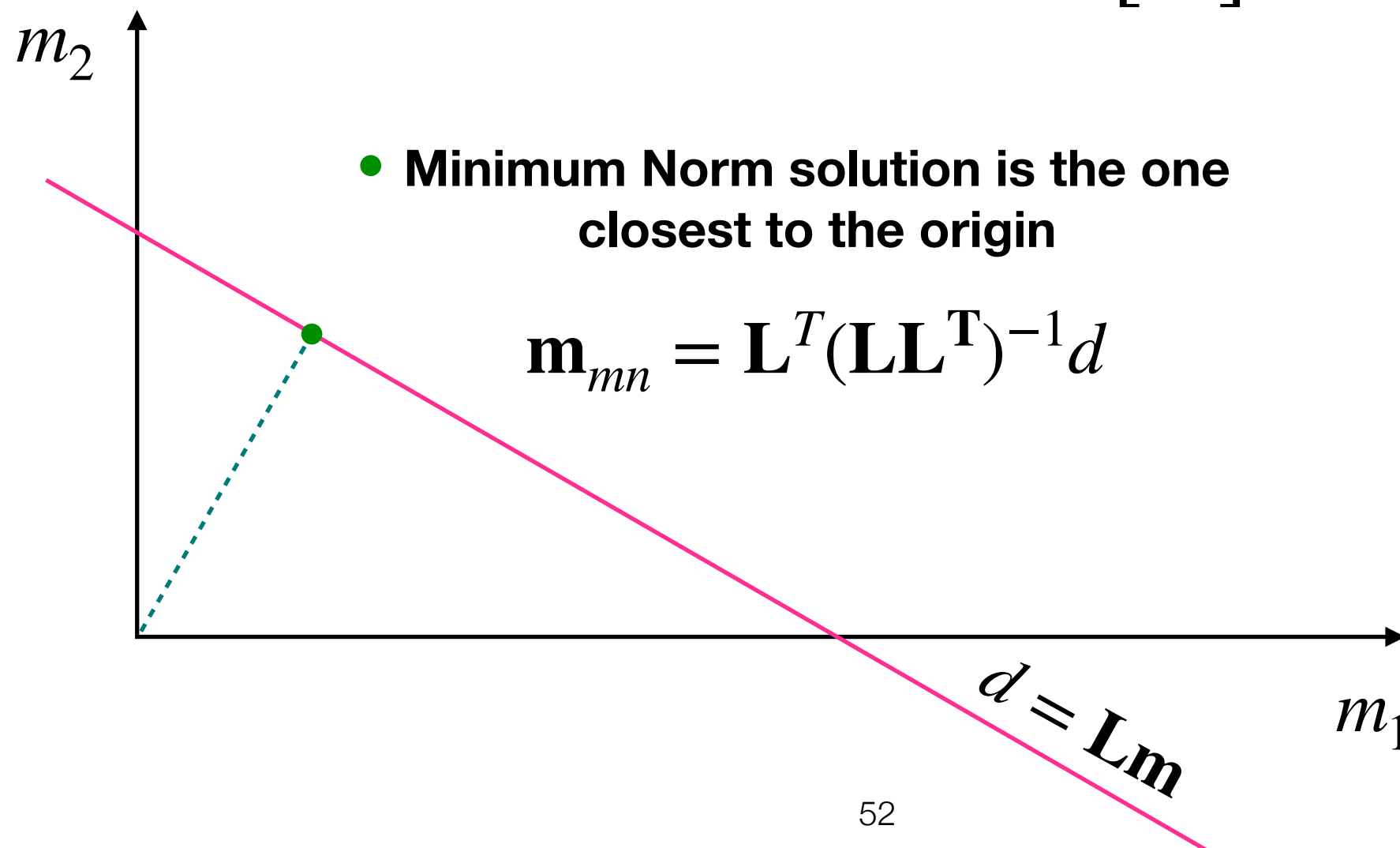
$$d = \mathbf{Lm} = \begin{bmatrix} L_{11} & L_{12} \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}$$



Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- $N=1, M=2$ (2 unknowns and 1 equation)

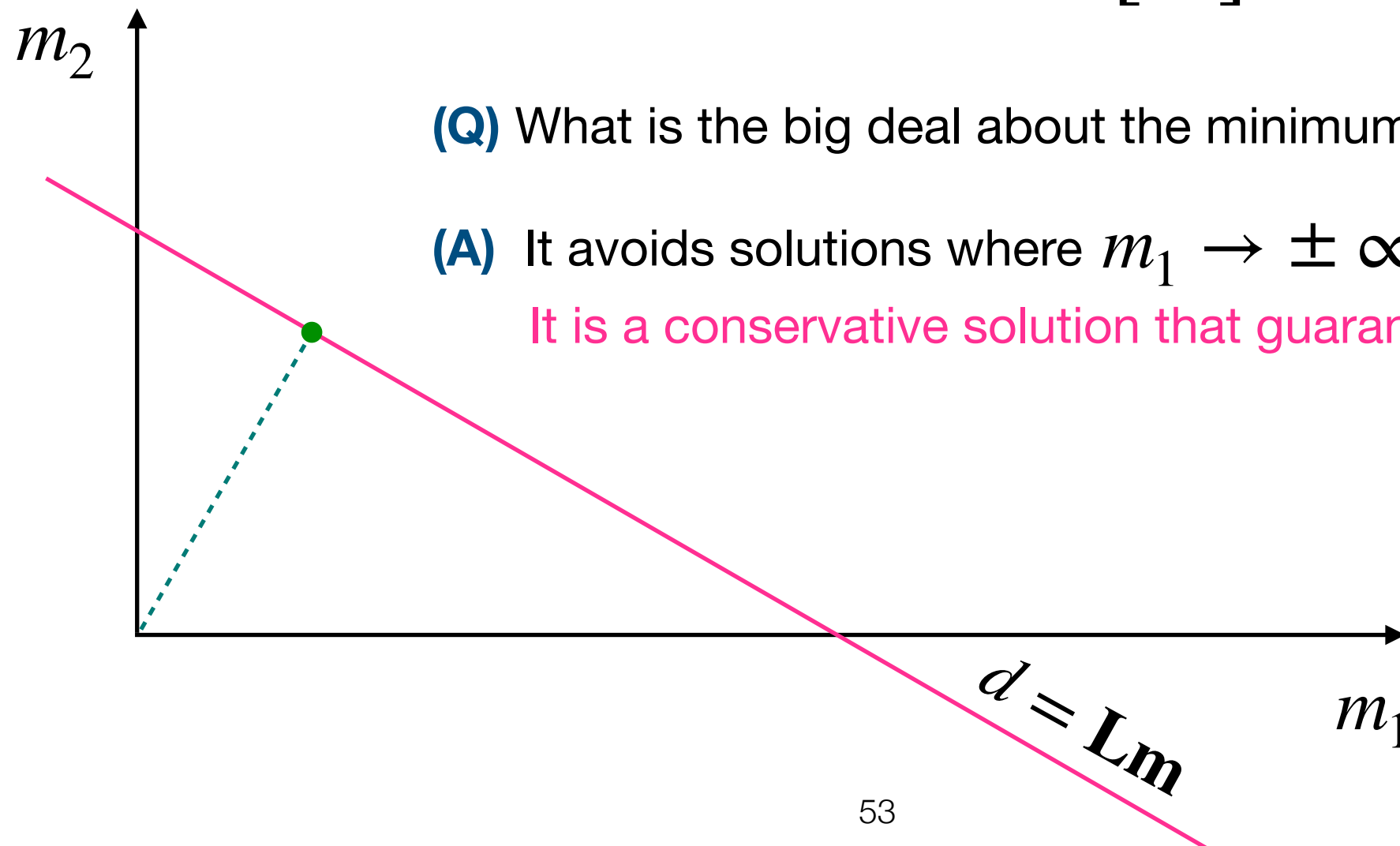
$$d = \mathbf{L}\mathbf{m} = \begin{bmatrix} L_{11} & L_{12} \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}$$



Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- $N=1, M=2$ (2 unknowns and 1 equation)

$$d = \mathbf{Lm} = \begin{bmatrix} L_{11} & L_{12} \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}$$



(Q) What is the big deal about the minimum norm solution?

(A) It avoids solutions where $m_1 \rightarrow \pm \infty, m_2 \rightarrow \pm \infty$

It is a conservative solution that guarantees stability

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

Problem

$$\mathbf{d} = \mathbf{L}\mathbf{m}$$

Minimum norm solution

$$\mathbf{m}_{mn} = \mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{d}$$

We can replace the data back into the solution $\mathbf{m}_{mn} = \mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{L}\mathbf{m}$

Model resolution matrix is not identity

$$\mathbf{m}_{mn} = \mathbf{R}_m\mathbf{m}$$

$$\mathbf{R}_m = \mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{L} \neq \mathbf{I}$$

So, unless the “true solution” is the minimum norm solution, we cannot recover the “true solution”...

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

Problem

$$\mathbf{d} = \mathbf{L}\mathbf{m}$$

Minimum norm solution

$$\mathbf{m}_{mn} = \mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{d}$$

Predicted data

$$\mathbf{d}^{pred} = \mathbf{L}\mathbf{m}_{mn} = \mathbf{L}\mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{d}$$

Data resolution matrix is the identity

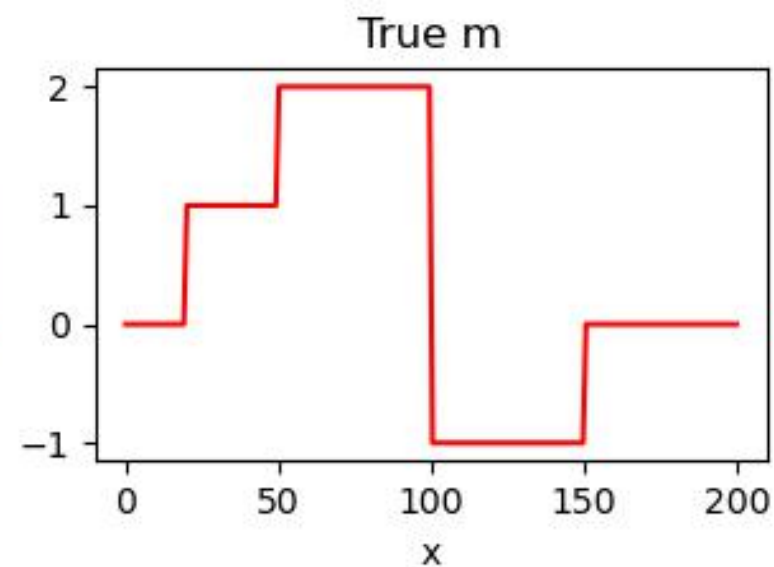
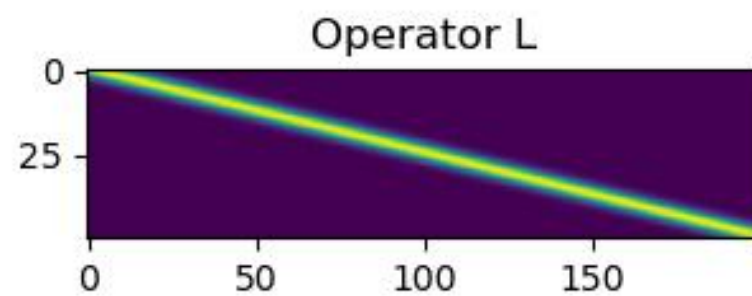
$$\mathbf{d}^{pred} = \mathbf{R}_d\mathbf{d} = \mathbf{d}$$

$$\mathbf{R}_d = \mathbf{L}\mathbf{L}^T(\mathbf{L}\mathbf{L}^T)^{-1} = \mathbf{I}$$

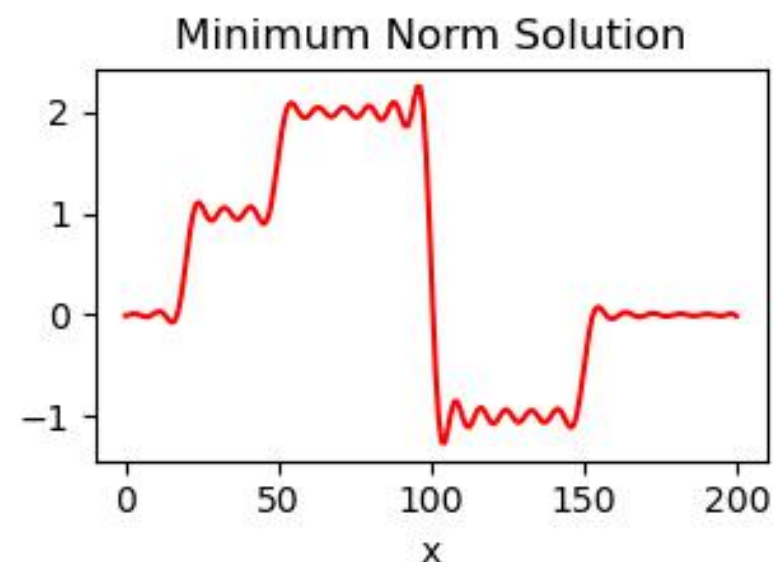
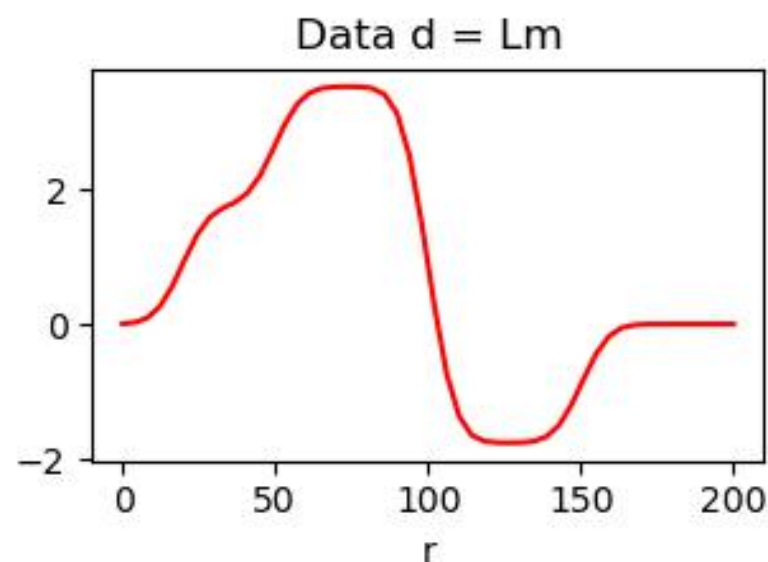
It is obvious that I can recover the data exactly because the data was used as an exact constraint when we constructed the minimum norm solution.

Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- Example $d(r_j) = \int_{x_{min}}^{x_{max}} A e^{-\alpha(r_j-x)^2} m(x) dx$

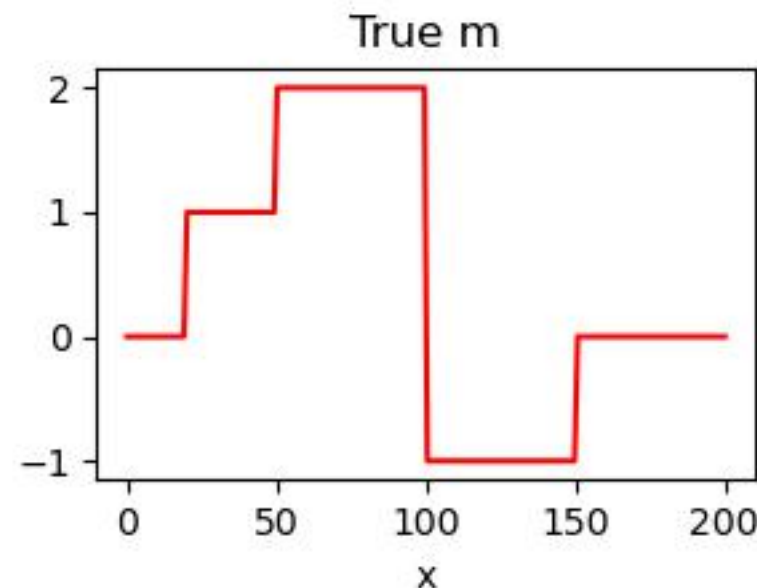
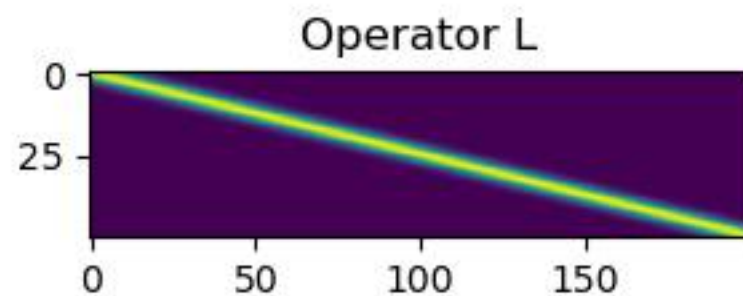


`demo_1.ipynb`

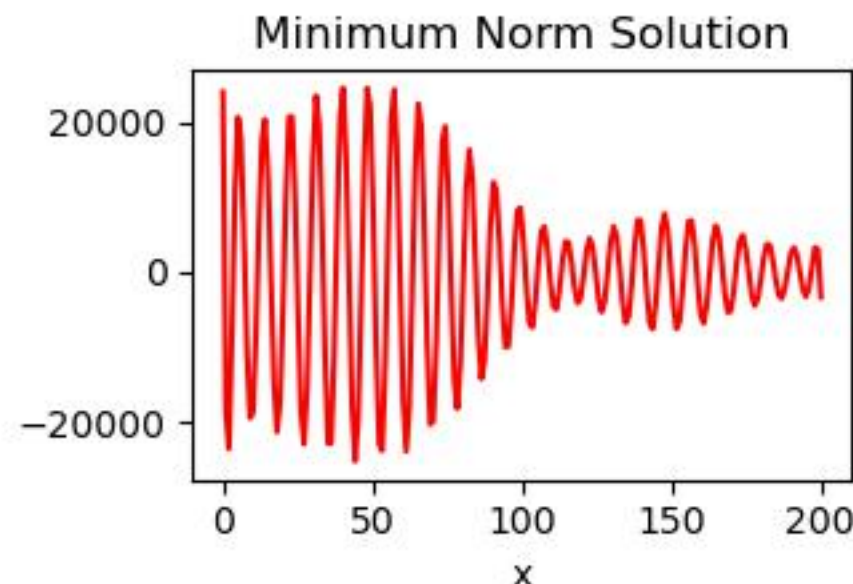
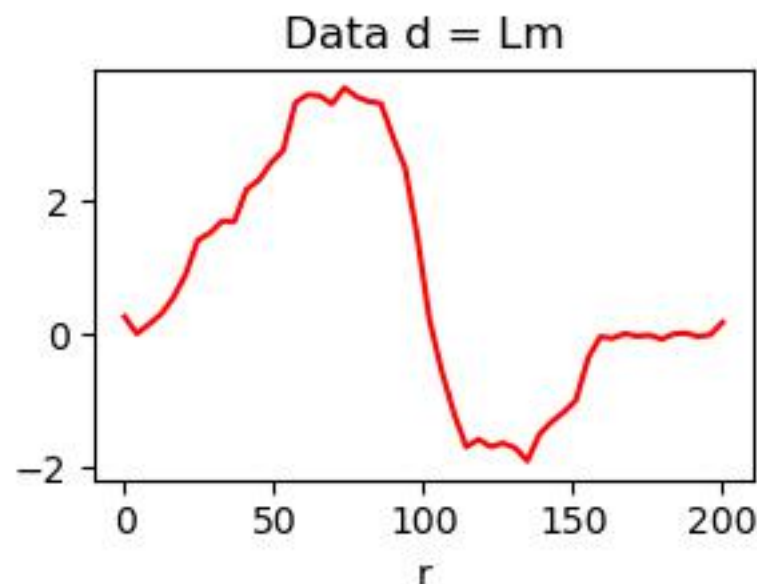


Discrete Linear Inverse Problem with Accurate Data: **Minimum Norm Solution**

- Example $d(r_j) = \int_x^{x_{max}} A e^{-\alpha(r_j-x)^2} m(x) dx + \sigma n(r_j)$



demo_2.ipynb



This is not working because I am fitting the data “exactly” and the data now contains noise!

Discrete Linear Inverse Problem with Accurate Data: **Weighted minimum Norm Solution**

- Among all possible solution select the one with minimum Weighted Euclidian norm. In other words,

$$\begin{aligned} &\textbf{Minimize } \|\mathbf{W}\mathbf{m}\|_2^2 \\ &\textbf{Subject to } \mathbf{d} = \mathbf{L}\mathbf{m} \end{aligned}$$

- The latter is solved using the Method of Lagrange Multipliers. We minimized an augmented cost function

$$J = \|\mathbf{W}\mathbf{m}\|_2^2 + \mathbf{b}^T(\mathbf{L}\mathbf{m} - \mathbf{d})$$

- $\mathbf{b} \in R^N$ is the vector of Lagrange Multipliers, one multiplier per observation

Discrete Linear Inverse Problem with Accurate Data: **Weighted minimum Norm Solution**

- To solve the problem we need to minimize the cost respect to \mathbf{m} and \mathbf{b}

$$J = \|\mathbf{W}\mathbf{m}\|_2^2 + \mathbf{b}^T(\mathbf{L}\mathbf{m} - \mathbf{d})$$

- Which leads to the so called weighted minimum norm solution (I will do it in the whiteboard)

$$\mathbf{Q} = (\mathbf{W}^T\mathbf{W})^{-1}$$
$$\mathbf{m}_{wmn} = \mathbf{Q}\mathbf{L}^T(\mathbf{L}\mathbf{Q}\mathbf{L}^T)^{-1}\mathbf{d}$$

Linear Inverse Problems with inaccurate data

- Including noise into the problem

$$\mathbf{d} = \mathbf{L}\mathbf{m} + \mathbf{e}$$

$$\mathbf{d} : N \times 1 \quad \mathbf{m} : M \times 1 \quad \mathbf{L} : N \times M$$

- \mathbf{e} represents “nice” Gaussian additive noise
- Given the vector of observed data \mathbf{d} , one wishes to estimate the vector of model parameters \mathbf{m}

Linear Inverse Problems with inaccurate data: **Regularized least-squares solution**

Cost function for least-squares problems

$$J = \|\mathbf{e}\|_2^2 = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2$$

The principle is simple, find \mathbf{m} that minimize the sum of the squares of the errors

$$\nabla J = 0 \rightarrow (\mathbf{L}^T \mathbf{L})\mathbf{m} = \mathbf{L}^T \mathbf{d}$$

To compute the solution we now have to invert a matrix.
Assume the matrix is invertible then

$$\mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{d} \quad M \ll N$$

Careful here: This is valid for overdetermined problems where the matrix to invert is full rank

Linear Inverse Problems with inaccurate data: **Regularized least-squares solution**

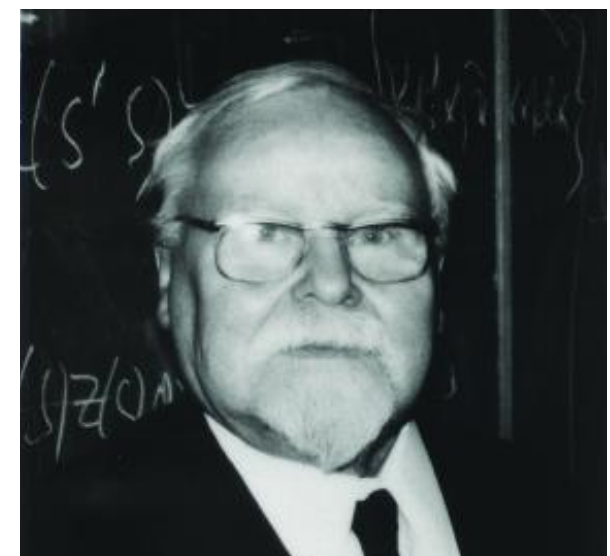
The latter is a **naive solution** because inverse problems are ill-posed and regularization is needed and, in general, they are not well-posed over-determined problems. In other words, one cannot safely compute: $(\mathbf{L}^T \mathbf{L})^{-1}$

- Either is non invertible
- Or the Matrix has a large condition number (It will amplify noise)
- Condition Number $\kappa = \frac{\lambda_{max}(\mathbf{L}^T \mathbf{L})}{\lambda_{min}(\mathbf{L}^T \mathbf{L})}$

Matrix computations
GH Golub, CF Van Loan, 2013 (4th edition)

Regularization (a solution to ill-conditioning)

- Main idea of regularization methods:
 - Take an ill-posed problem and turn into a well-posed problem by introducing constraints that lead to a stable solution. The solution often depends on a **trade-off** parameter. Therefore, regularization methods create a family of solutions with different properties (e.g. smoothness). One can modify the **solution** by changing the tradeoff parameter
 - Often called Tikhonov regularization



A. N. Tikhonov (1906-1993)

Regularization

- Replace cost function to minimize

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 \longleftarrow \text{Cost with no regularization term}$$

- by

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W}\mathbf{m}\|_2^2$$

- μ Is the infamous trade-off parameter or regularization parameter. Why infamous?
- \mathbf{W} is a matrix/operator of weights

Regularization

- Evaluate the solution, as usual, by minimizing a cost function

$$\begin{aligned}\mathbf{m}_{sol} &= \underset{\mathbf{m}}{\operatorname{argmin}} J \\ &= \underset{\mathbf{m}}{\operatorname{argmin}} \{ \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W} \mathbf{m}\|_2^2 \}\end{aligned}$$

- Taking derivatives and equating them to zero

$$\nabla J = 0 \rightarrow \mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{R})^{-1} \mathbf{L}^T \mathbf{d}$$

$$\mathbf{R} = \mathbf{W}^T \mathbf{W}$$

Regularization

- Anatomy of the cost function

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{W}\mathbf{m}\|_2^2}_{\text{Model Norm}}$$

$$J = \underbrace{\|\mathbf{e}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{u}\|_2^2}_{\text{Model Norm}}$$

- To minimize J is equivalent to simultaneously minimize \mathbf{e} and \mathbf{u}

Regularization

- To minimize J is equivalent to simultaneously minimize \mathbf{e} and \mathbf{u}

$$\mathbf{e} = \mathbf{d} - \mathbf{L}\mathbf{m} \approx \mathbf{0} \quad \text{Minimize the residuals}$$

$$\mathbf{u} = \mathbf{W}\mathbf{m} \approx \mathbf{0} \quad \text{Minimize the bad features of } \mathbf{m}$$

- \mathbf{W} is a high-pass operator therefore it penalizes roughness. The vector \mathbf{u} represents amplified bad features of \mathbf{m}
- Examples of \mathbf{W} are first and second order derivatives

Regularization with **first order** derivative smoothing

$$\mathbf{W} = \mathbf{D}_1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{pmatrix}$$

$$\frac{\partial f(x)}{\partial x} \leftrightarrow ikF(k)$$

The derivative operator is high pass because it amplifies high frequencies

Regularization with **second order** derivative smoothing

$$\mathbf{W} = \mathbf{D}_2 = \begin{pmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{pmatrix}$$

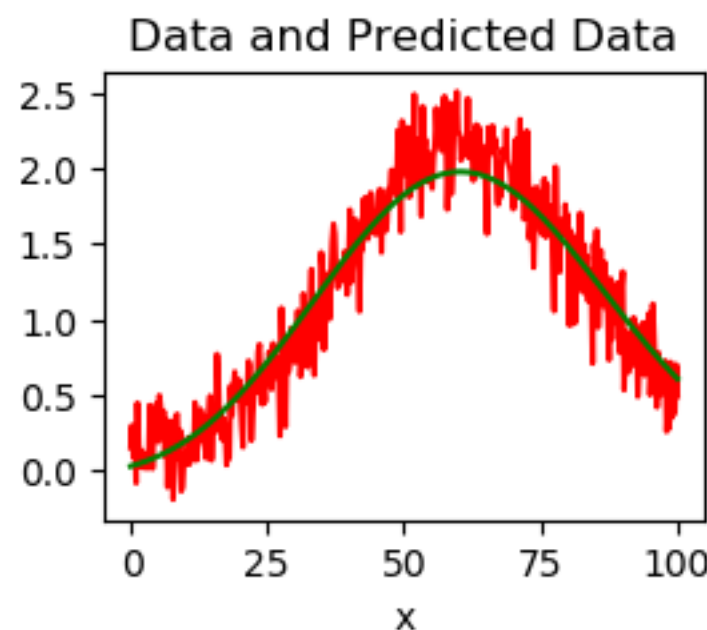
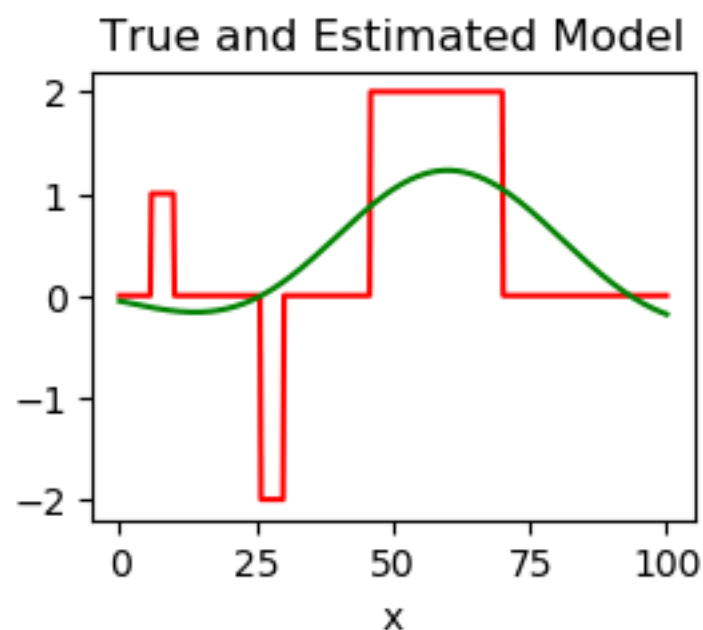
$$\frac{\partial^2 f(x)}{\partial x^2} \leftrightarrow -k^2 F(k) \quad \text{Second Derivative is also a high-pass operator}$$

Damped LS solution is when $W=I$

$$\mathbf{m}_{sol} = \underset{\mathbf{m}}{\operatorname{argmin}} J$$

$$= \underset{\mathbf{m}}{\operatorname{argmin}} \{ \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2 \}$$

$$\frac{\partial J}{\partial \mathbf{m}} = 0 \rightarrow \mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d}$$



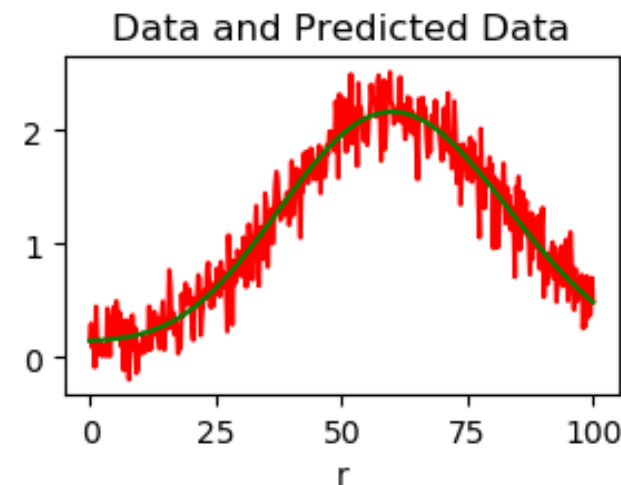
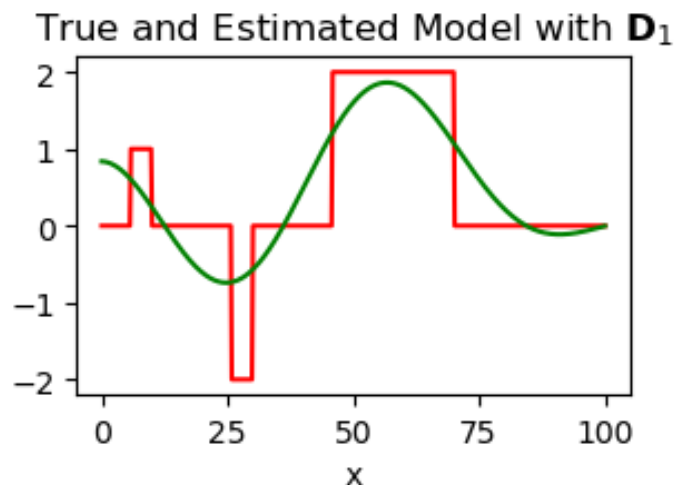
Red: True Model | Observed data
Green: Estimated Model | Predicted data

Example: Solution with smoothing

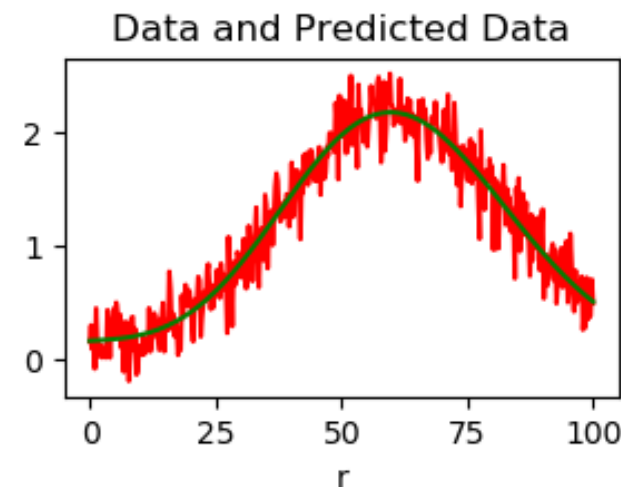
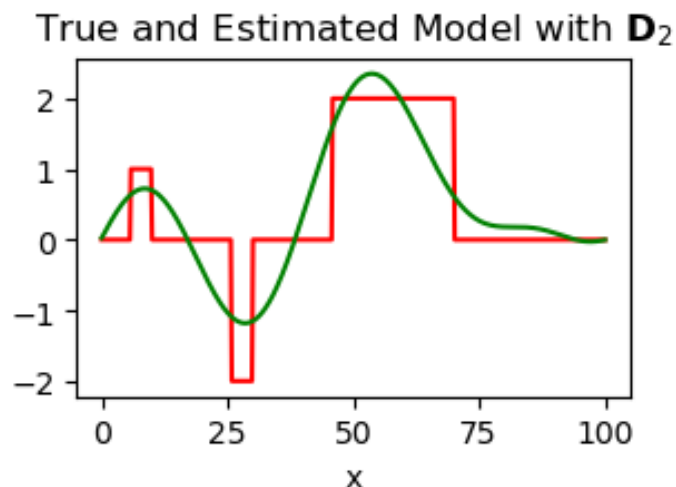
$$\mathbf{m}_{sol} = \underset{\mathbf{m}}{\operatorname{argmin}} \{ \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W}\mathbf{m}\|_2^2 \}$$

$$\frac{\partial J}{\partial \mathbf{m}} = 0 \rightarrow \mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{W}^T \mathbf{W})^{-1} \mathbf{L}^T \mathbf{d}$$

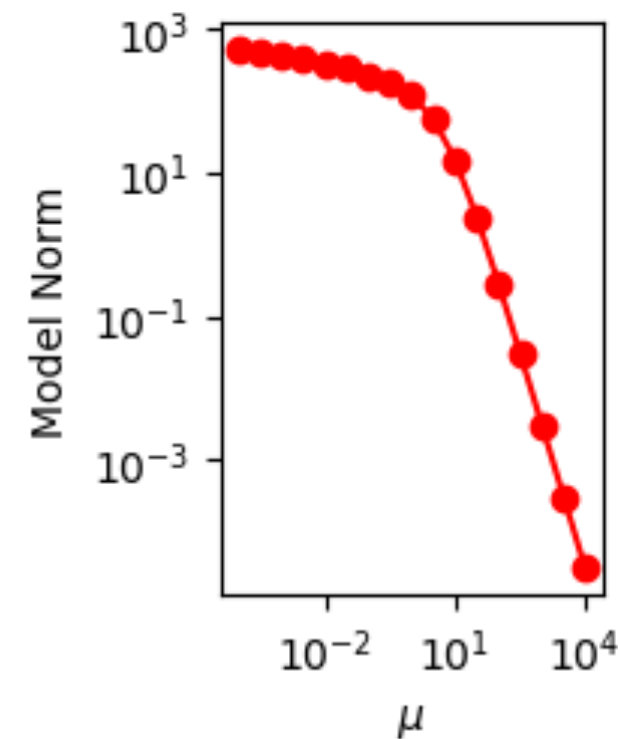
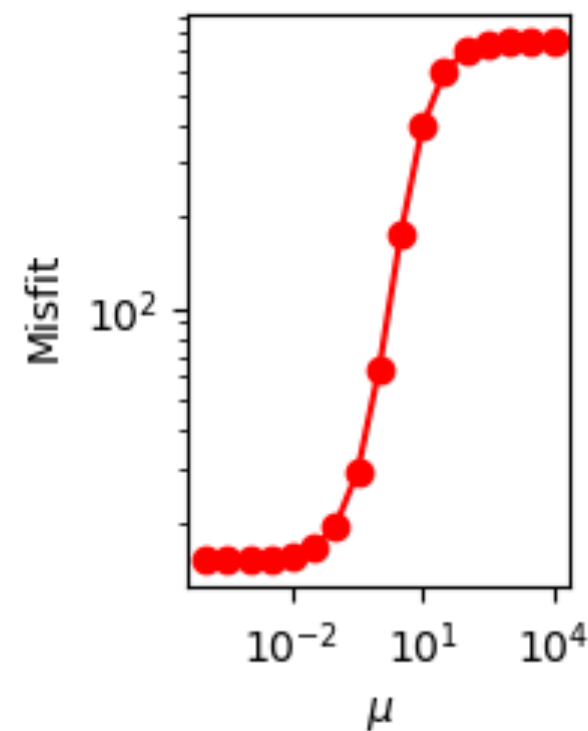
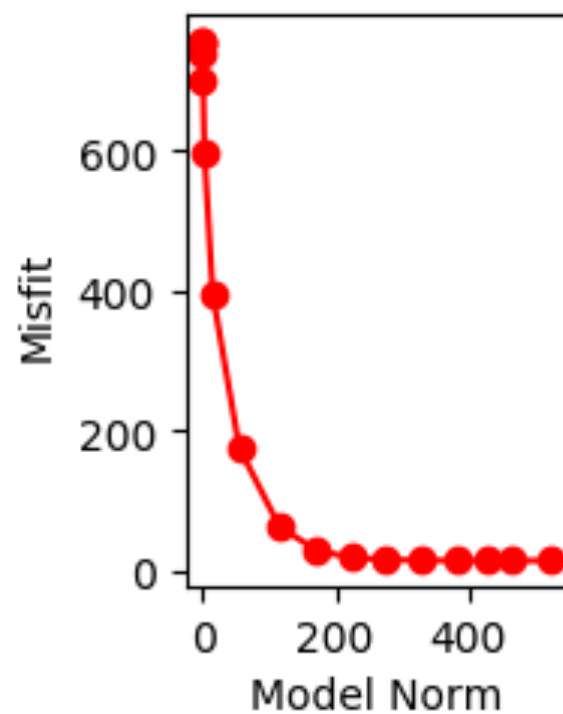
$$\mathbf{W} = \mathbf{D}_1$$



$$\mathbf{W} = \mathbf{D}_2$$



Trade-off curves (for Damped least-squares)



Underdetermined Problem with inaccurate data

Minimum Norm Solution

$$J = \text{Misfit} + \mu \text{ Model Norm}$$

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2$$

Make residuals small = Find a model that fits the data



Make model parameters small = Stable solution that does not blow up

μ : Trade-off parameter of the inverse problem

Underdetermined Problem with inaccurate data

Minimum Norm Solution

- So we solve the problem

$$\mathbf{d} = \mathbf{L}\mathbf{m} + \mathbf{n}$$

by minimizing the cost $J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu\|\mathbf{m}\|_2^2$

- We already discussed how one can compute derivatives of the scalar function J with respect to a vector.
- Condition for minimum

$$\frac{dJ}{d\mathbf{m}} = 0 \rightarrow (\mathbf{L}^T\mathbf{L} + \mu\mathbf{I})\mathbf{m} = \mathbf{L}^T\mathbf{d}$$

Underdetermined Problem with inaccurate data

Minimum Norm Solution

$$\text{Cost } J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2$$

- Solution $\mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d}$
- Predicted data $\mathbf{d}_{pred} = \mathbf{L}\mathbf{m}_{sol} = \mathbf{L}(\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d}$
- Predicted residuals $\mathbf{r} = \mathbf{d}_{pred} - \mathbf{d}$
- You can also show the following important identity:

$$\mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d} = \mathbf{L}^T (\mathbf{L} \mathbf{L}^T + \mu \mathbf{I})^{-1} \mathbf{d}$$

Underdetermined Problem with inaccurate data

Minimum Norm Solution

$$\mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d} = \mathbf{L}^T (\mathbf{L} \mathbf{L}^T + \mu \mathbf{I})^{-1} \mathbf{d}$$

1

2

1 looks like the least-squares solution for overdetermined problems, so it is often called the regularized least-squares solution. Also called damped least-squares, ridge regression solution, least-squares solution with zero-order quadratic regularization, Tikhonov solution, etc.

2 looks like the minimum norm solution for exact data but now there is an additional constant diagonal term.

The method we described is also called Tikhonov regularization, so the solution given by either 1) or 2) is also the Tikhonov solution.

Underdetermined Problem with inaccurate data

Minimum Norm Solution

$$\mathbf{m}_{sol} = (\mathbf{L}^T \mathbf{L} + \mu \mathbf{I})^{-1} \mathbf{L}^T \mathbf{d} = \mathbf{L}^T (\mathbf{L} \mathbf{L}^T + \mu \mathbf{I})^{-1} \mathbf{d}$$

Where does the name damped least-squares solution come from? As one increases μ , more damping is introduced. In other words, the solution becomes smaller and less oscillatory!

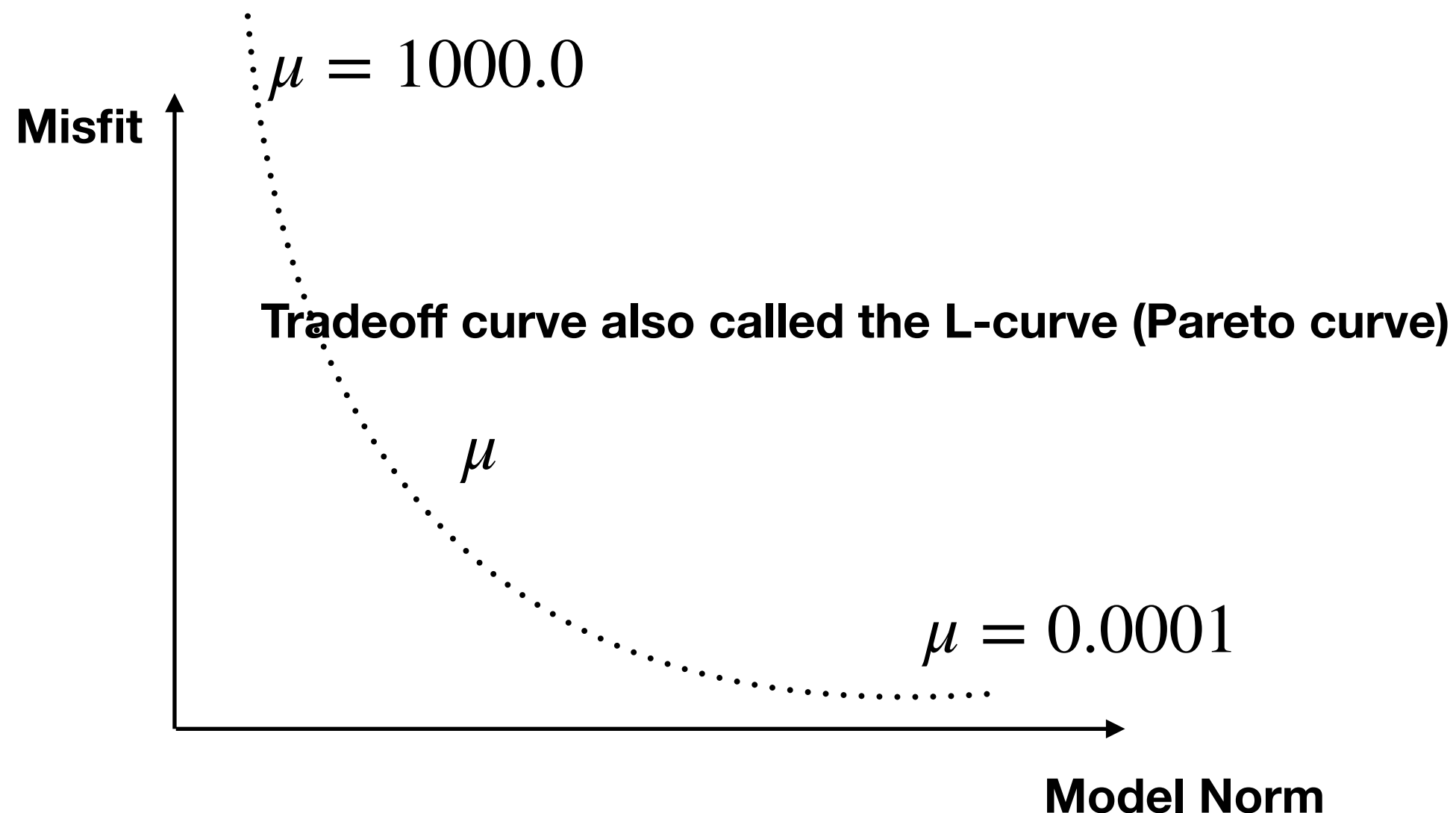
The tradeoff parameter μ is sometimes also call the damping parameter.

The tradeoff parameter must be provided and it defines a family of solutions. I should have written $\mathbf{m}_{sol}(\mu)$ to indicate the dependency of the solution on the tradeoff parameter.

Underdetermined Problem with inaccurate data

Minimum Norm Solution

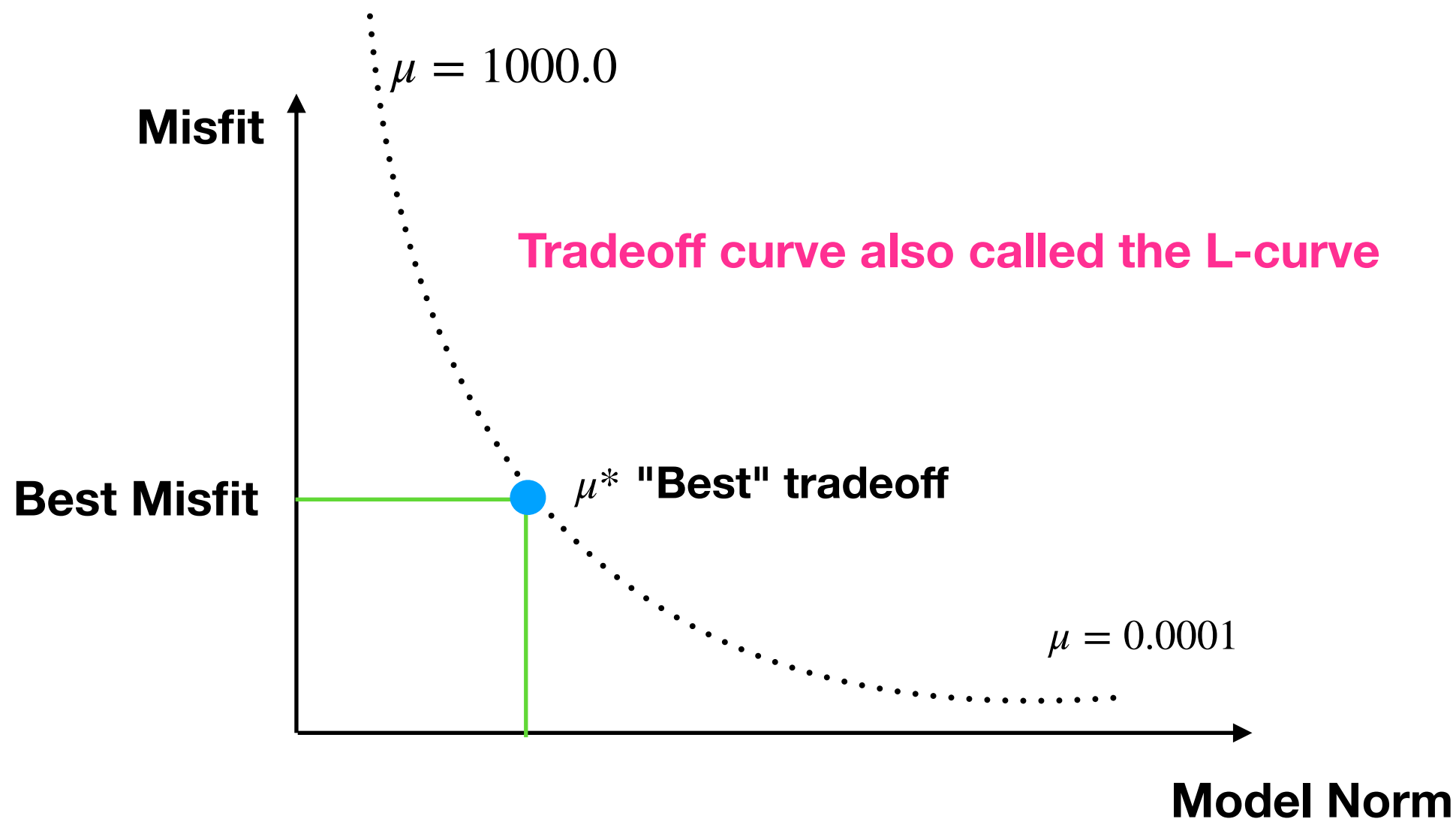
$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2 = \mathbf{Misfit} + \mu \mathbf{Model\ Norm}$$



Underdetermined Problem with inaccurate data

Minimum Norm Solution

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2 = \mathbf{Misfit} + \mu \mathbf{Model\ Norm}$$



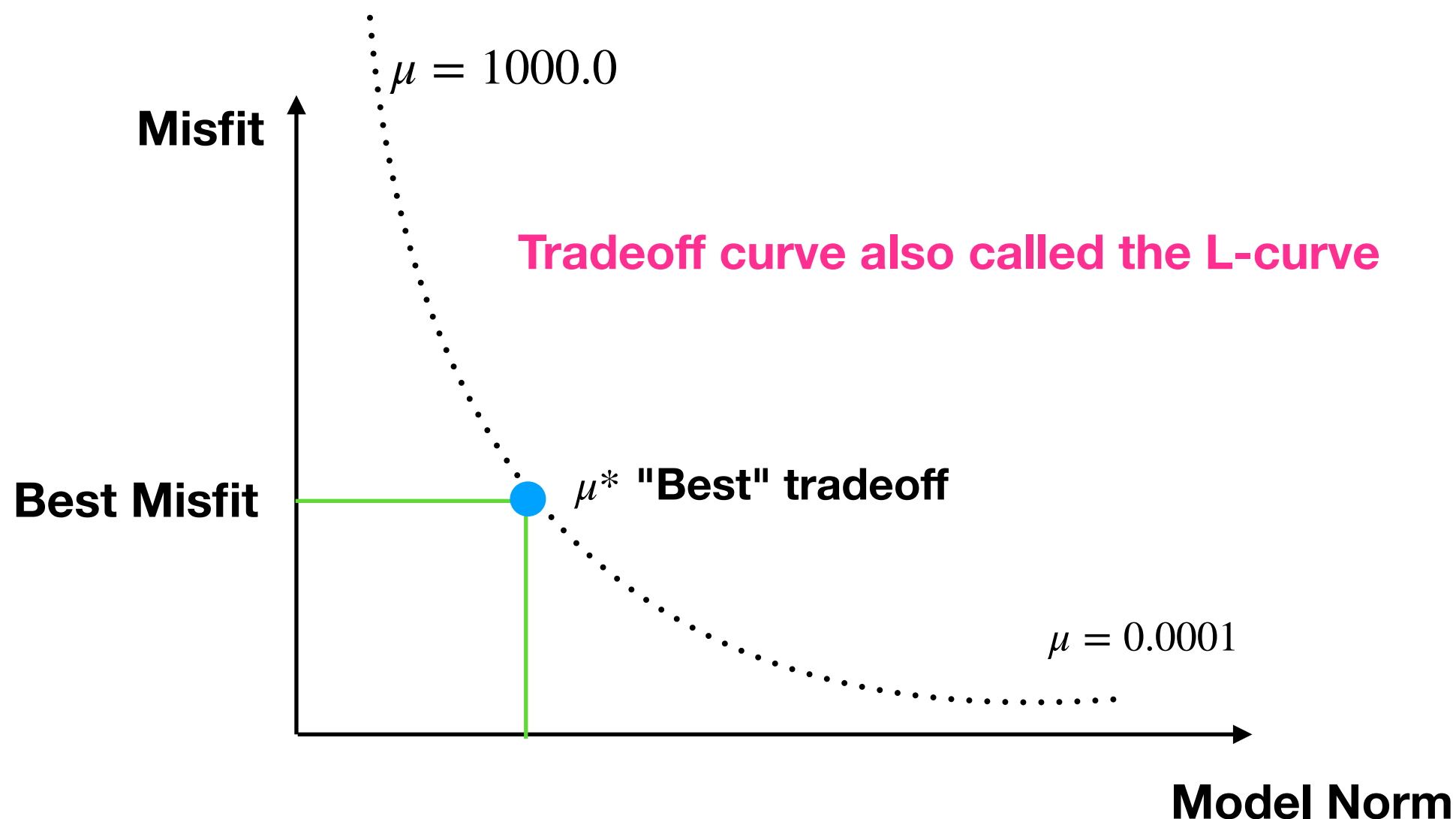
Underdetermined Problem with inaccurate data

Minimum Norm Solution

Recall original problem: $\mathbf{d} = \mathbf{L}\mathbf{m} + \mathbf{n}$.

Assume you know the variance of the noise: σ^2

$$\text{Best misfit} = \|\mathbf{d} - \mathbf{L}\mathbf{m}_{sol}(\mu^*)\|_2^2 = N \times \sigma^2$$



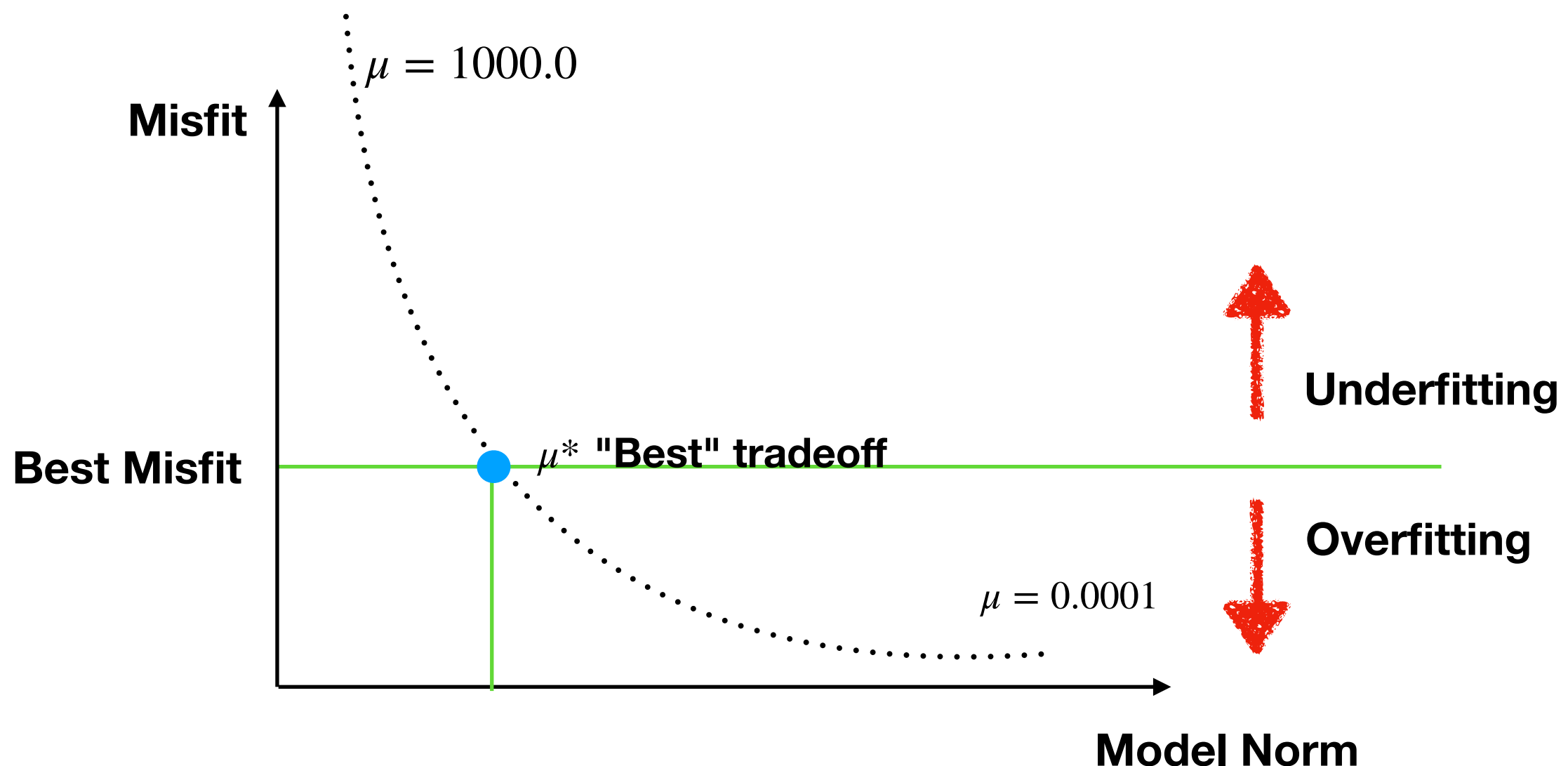
Underdetermined Problem with inaccurate data

Minimum Norm Solution

Recall original problem: $\mathbf{d} = \mathbf{L}\mathbf{m} + \mathbf{n}$.

Assume you know the variance of the noise: σ^2

$$\text{Best misfit} = \|\mathbf{d} - \mathbf{L}\mathbf{m}_{\text{sol}}(\mu^*)\|_2^2 = N \times \sigma^2$$

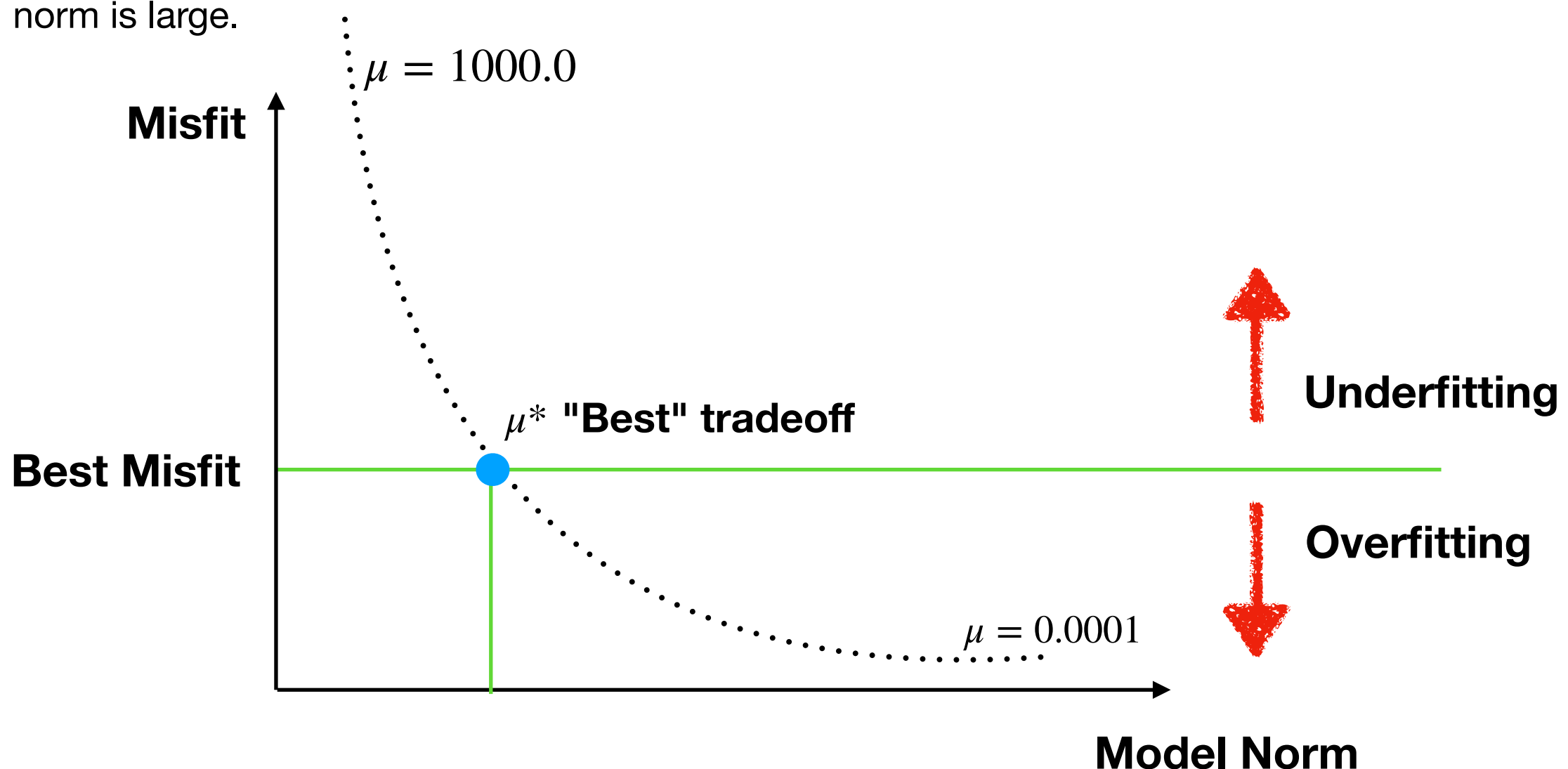


Underdetermined Problem with inaccurate data

Minimum Norm Solution

Underfitting: Observations are not properly fit. You are overestimating the amount of noise in the data. Solution norm decreases. Residuals are too large. There is probably some structure in the residuals.

Overfitting: Residuals become too small and therefore your estimated model produces data that fits the noise. As a consequence of the latter, the model can become unstable and the model norm is large.



Underdetermined Problem with inaccurate data

Weighted Minimum Norm Solution

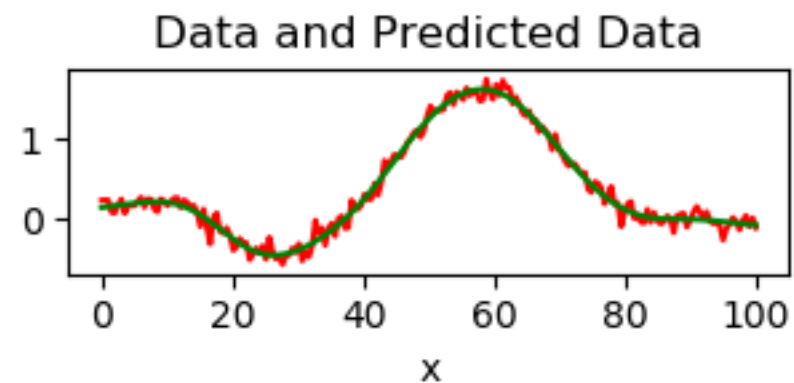
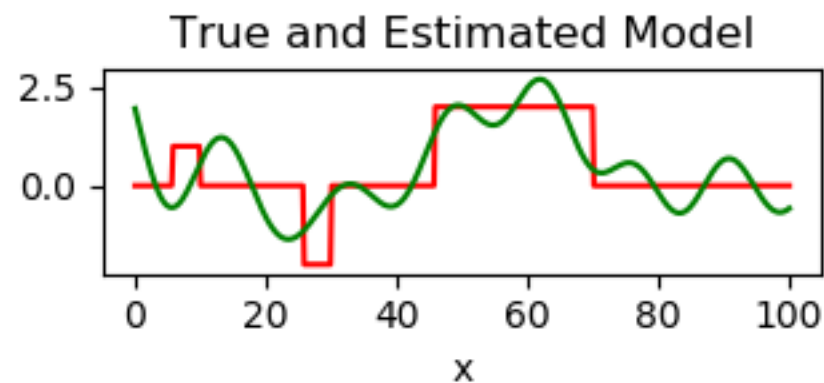
$$\text{Cost } J = \|\mathbf{d}_{obs} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W}\mathbf{m}\|_2^2$$

- Solution $\mathbf{m}_{sol} = (\mathbf{L}^T\mathbf{L} + \mu\mathbf{W}^T\mathbf{W})^{-1}\mathbf{L}^T\mathbf{d}$
- Predicted data $\mathbf{d}_{pred} = \mathbf{L}\mathbf{m}_{sol} = \mathbf{L}(\mathbf{L}^T\mathbf{L} + \mu\mathbf{W}^T\mathbf{W})^{-1}\mathbf{L}^T\mathbf{d}_{obs}$
- Predicted residuals $\mathbf{r} = \mathbf{d}_{pred} - \mathbf{d}$
- You can also show the following important identity:

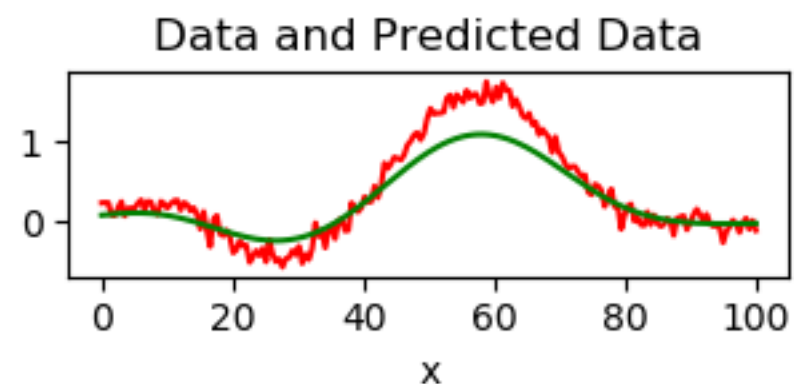
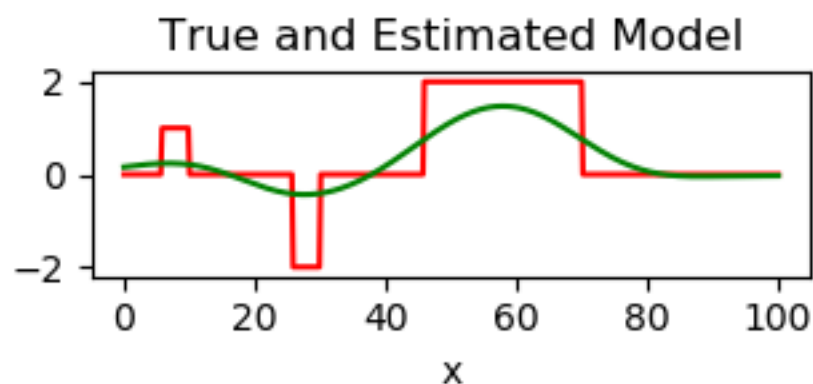
$$\mathbf{m}_{sol} = (\mathbf{L}^T\mathbf{L} + \mu\mathbf{Q})^{-1}\mathbf{L}^T\mathbf{d} = \mathbf{Q}^{-1}\mathbf{L}^T(\mathbf{L}\mathbf{Q}^{-1}\mathbf{L}^T + \mu\mathbf{I})^{-1}\mathbf{d}$$

$$\mathbf{Q} = \mathbf{W}^T\mathbf{W}$$

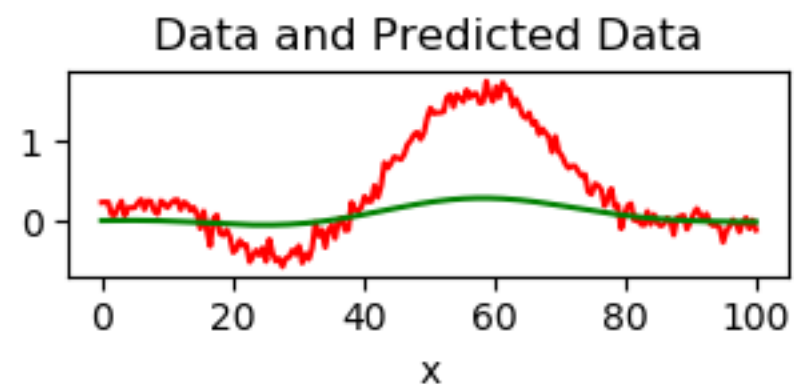
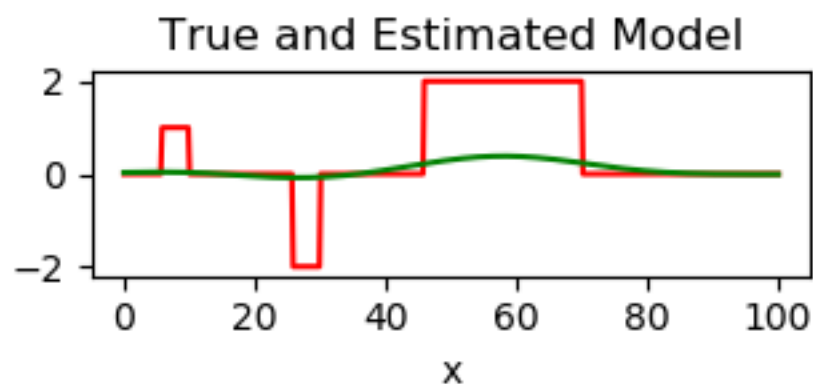
Trade-off curves (for Damped least-squares)



$$\mu = 0.0001$$



$$\mu = 0.1$$



$$\mu = 1$$

Changing the norm

Edge preserving regularization (EPR)

- Avoid smoothing to preserve edges
- We adopt the ell-1 norm of the first order derivative of model parameters
- Make the derivative of the model sparse

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{D}_1 \mathbf{m}\|_1}_{\text{Model Norm}}$$

We ask $\mathbf{u} = \mathbf{D}_1 \mathbf{m}$ to be sparse

The L_1 and the L_2 norms

$$\mathbf{u} : N \times 1 \longrightarrow \|\mathbf{u}\|_2^2 = \mathbf{u}^H \mathbf{u} = \sum_{i=1}^N u_i u_i^* = \sum_{i=1}^N |u_i|^2$$

$$\mathbf{u} : N \times 1 \longrightarrow \|\mathbf{u}\|_1 = \sum_{i=1}^N |u_i|$$

Edge preserving regularization (EPR)

- Quadratic leads to close form solution (linear system of equations)

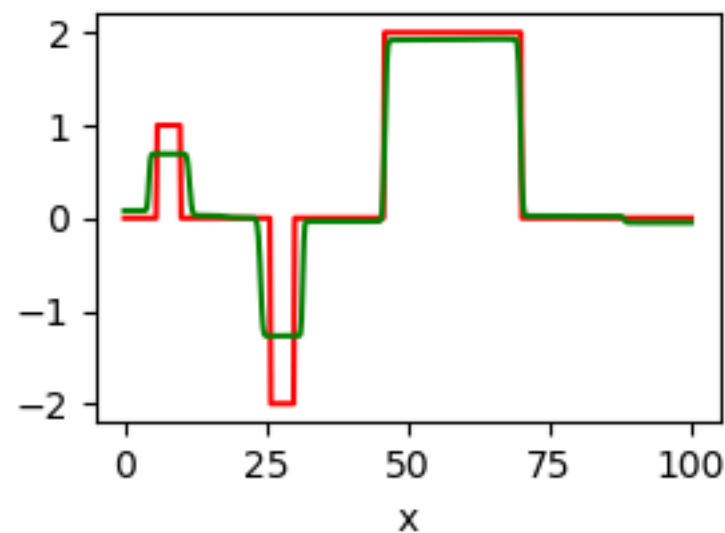
$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{W}\mathbf{m}\|_2^2}_{\text{Model Norm}}$$

- **EPR:** Non-quadratic regularization leads to non-linear solution that tries to recover edges

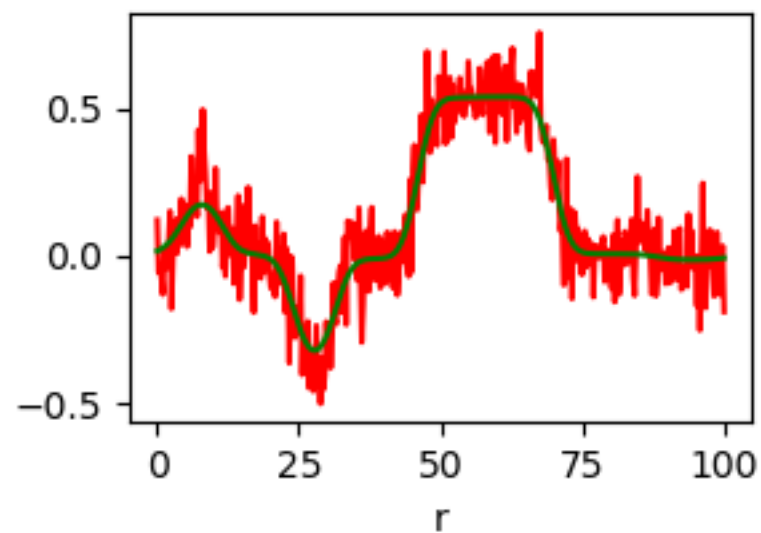
$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{D}_1\mathbf{m}\|_1}_{\text{Model Norm}}$$

Edge preserving regularization (EPR)

True and Estimated Model with EPR

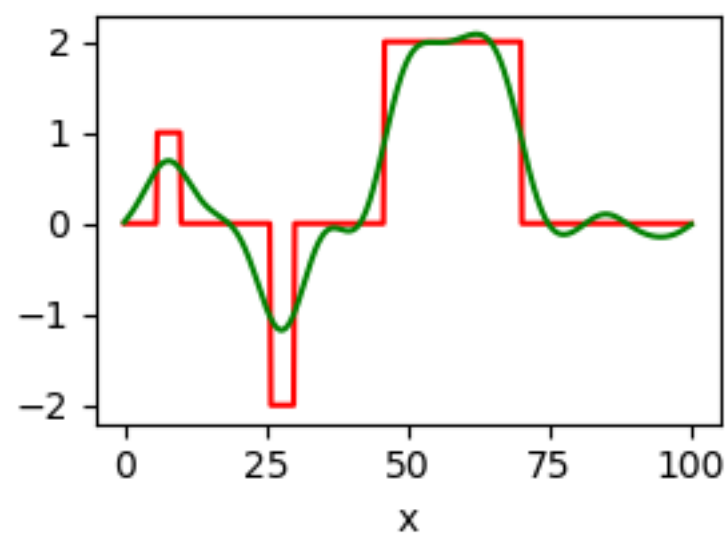


Data and Predicted Data

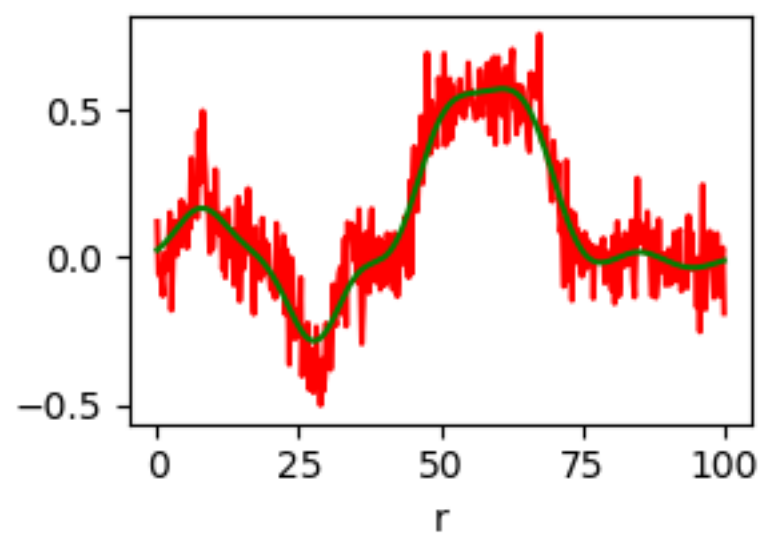


EPR
(non-quadratic)

True and Estimated Model with \mathbf{D}_2



Data and Predicted Data



Second Order Derivative Reg.
(quadratic)

Edge preserving regularization (EPR)

- Non-quadratic regularization leads to non-linear solution

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{D}_1 \mathbf{m}\|_1}_{\text{Model Norm}}$$

- Solution

$$\nabla J = 0 \rightarrow (\mathbf{L}'\mathbf{L} + \mu\mathbf{D}_1' \mathbf{Q} \mathbf{D}_1)\mathbf{m} = \mathbf{L}'\mathbf{d}$$

- Where $\mathbf{v} = \mathbf{D}_1 \mathbf{m}$ $Q_{ii} = \frac{1}{|v_i|}$

- To avoid division by zero $Q_{ii} = \frac{1}{\epsilon + |v_i|}$

Edge preserving regularization (EPR)

- Iterative re-weighted least-squares

$$\mathbf{m}^1 = \mathbf{m}_{initial}$$

For $k=1$ **until convergence**

$$\mathbf{v} = \mathbf{D}_1 \mathbf{m}$$

$$Q_{ii}^k = \frac{1}{\epsilon + |v_i^k|}$$

$$\mathbf{m}^{k+1} = (\mathbf{L}'\mathbf{L} + \mu\mathbf{D}_1' \mathbf{Q}^k \mathbf{D}_1)^{-1} \mathbf{L}'\mathbf{d}$$

End

Connection to sparsity

- This cost function generates a sparse solution

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{m}\|_1}_{\text{Model Norm}}$$

- Often used for Deconvolution
- Pre-stack data Reconstruction (*Liu and Sacchi, 2004, Geophysics; Hermann, 2010, Geophysics*)
- Radon Transforms, etc etc etc (*Sacchi & Ulrych, 1995, Geophysics*)
- AVO Inversion (*Alemie and Sacchi, 2011, Geophysics*)

Connection to sparsity

- Make \mathbf{m} sparse:

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{m}\|_1}_{\text{Model Norm}}$$

- Make the derivative of \mathbf{m} sparse = Make \mathbf{m} blocky

$$J = \underbrace{\|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2}_{\text{Error Norm}} + \mu \underbrace{\|\mathbf{D}_1 \mathbf{m}\|_1}_{\text{Model Norm}}$$

Sparsity (more to follow...)

- **IRLS is the simplest solver one can imagine**
- Many new solvers in recent years
 - ISTA, FISTA, SALSA, SPG-L1, ADMM, L1-Magic, etc etc etc.. show **Notes in Latex**
 - I usually use IRLS or FISTA
 - ISTA/FISTA: Only need to know how to apply \mathbf{L} and \mathbf{L}' (*on-the-flight*)

PROBLEMS IN IMPLICIT FORM

Forward and Adjoint Operators

- Two special operations (or operators?)

Forward : $\mathbf{d} = \mathbf{L}\mathbf{m}$

Conjugate transpose or adjoint : $\tilde{\mathbf{m}} = \mathbf{L}'\mathbf{d}$

- \mathbf{L} is a linear operator (Forward modelling operator).
- Why are them special?
 - Iterative solvers for large inverse problems only need to know how to evaluate $\mathbf{L}[\cdot]$ and $\mathbf{L}'[\cdot]$
 - I usually interpret operators as matrices. In reality, operators are codes applied *on the flight*
 - *We will see these operators everywhere today*

Forward and Adjoint Operators

- **For matrices**

Forward : L

Conjugate transpose or adjoint : $L' = L^T$

- Examples of parts (L, L')
 - Demigration, Migration
 - Radon Modeling \mathcal{R} , Radon adjoint \mathcal{R}'
 - Sampling, Inserting zeros
 - Summation, Distribution or spraying
 - Fourier Synthesis \mathcal{F} , Fourier Analysis \mathcal{F}^H

Steepest descent method

- Compute the gradient and iterate downhill until convergence
- Let's see the special role of \mathbf{L} and \mathbf{L}' in steepest descent optimization (or in any iterative optimization algorithm that only requires \mathbf{L} and \mathbf{L}')
- Simplify problem by considering minimization of quadratic cost J

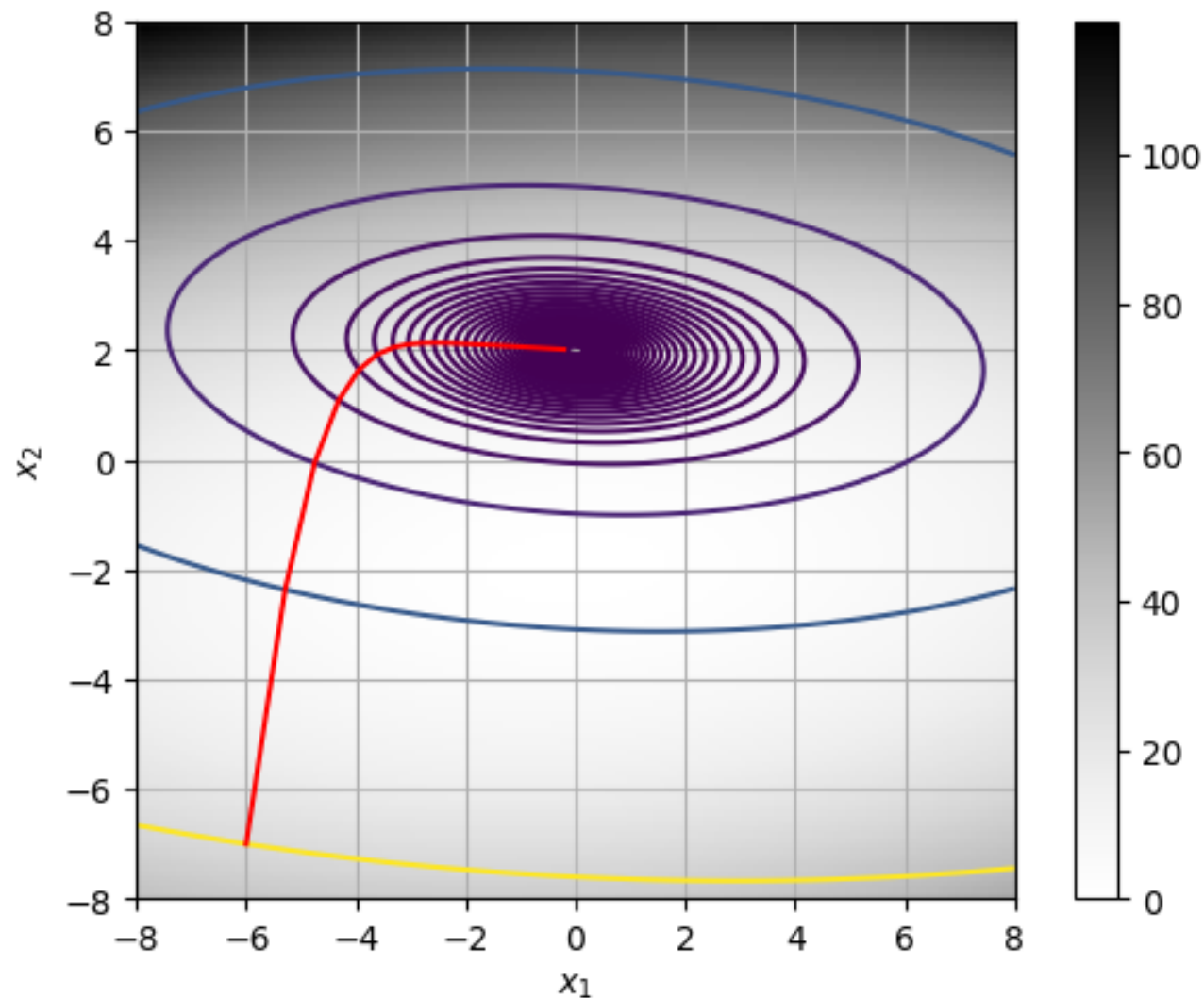
$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2$$

Steepest descent method

- Cost $J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2$
- Gradient $\mathbf{g} = \nabla J = \mathbf{L}'(\mathbf{L}\mathbf{m} - \mathbf{d})$
- Update $\mathbf{m}^{k+1} = \mathbf{m}^k - \alpha \mathbf{g}^k$
- Interesting, one can minimize J with a simple rule that does not involve inverting matrices. This is great because we can replace matrices by linear operators!!!!
- Aren't you excited ?

Steepest descent (SD) method

- Minimization of a quadratic cost by SD



$x_{\text{init}} = [-6.0, -7.9]$

$x_{\text{final}} = [-0.009, 2.005]$

$x_{\text{true}} = [0.0, 2.0]$

$K = 50$ iterations at fixed step size

Steepest descent method

SD iterations:

```
Do until convergence
     $\mathbf{m}^{k+1} = \mathbf{m}^k - \alpha(\mathbf{L}'(\mathbf{Lm}^k - \mathbf{d}))$ 
End Do
```

The above code can also be written as follows

```
Do until convergence
     $\mathbf{r}^k = (\mathbf{Lm}^k - \mathbf{d})$ 
     $\mathbf{m}^{k+1} = \mathbf{m}^k - \alpha \mathbf{L}'(\mathbf{r}^k)$ 
End Do
```

Steepest descent method: Matrices are replaced by Linear Operators packed into Functions or Subroutines

Do until convergence

$$\mathbf{r}^k = (\mathbf{L}\mathbf{m}^k - \mathbf{d})$$

$$\mathbf{m}^{k+1} = \mathbf{m}^k - \alpha \mathbf{L}'(\mathbf{r}^k)$$

End Do

Do until convergence

$$\mathbf{r}^k = \mathbf{Do_It}[\mathbf{m}^k, flag = f] - \mathbf{d}$$

$$\mathbf{m}^{k+1} = \mathbf{m}^k - \alpha \mathbf{Do_It}[\mathbf{r}^k, flag = a]$$

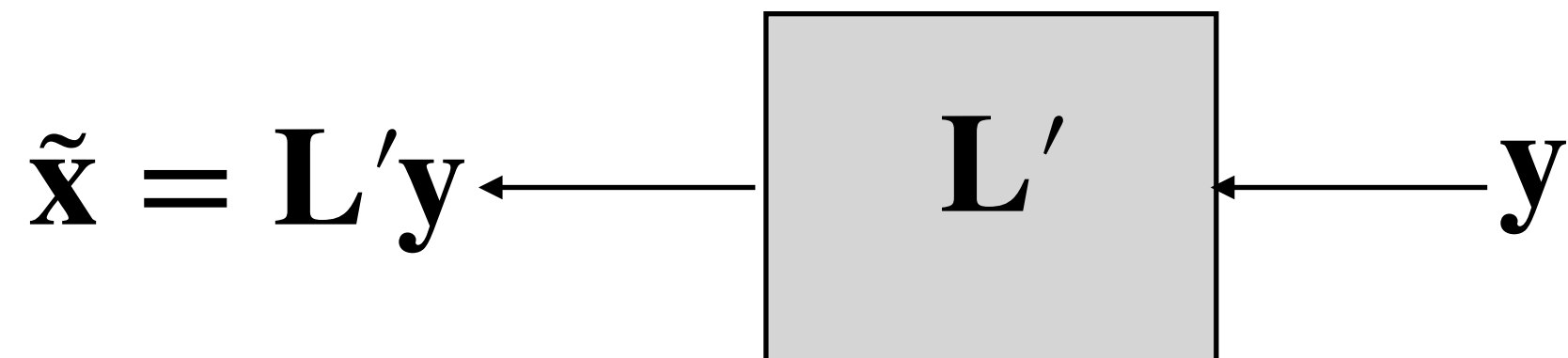
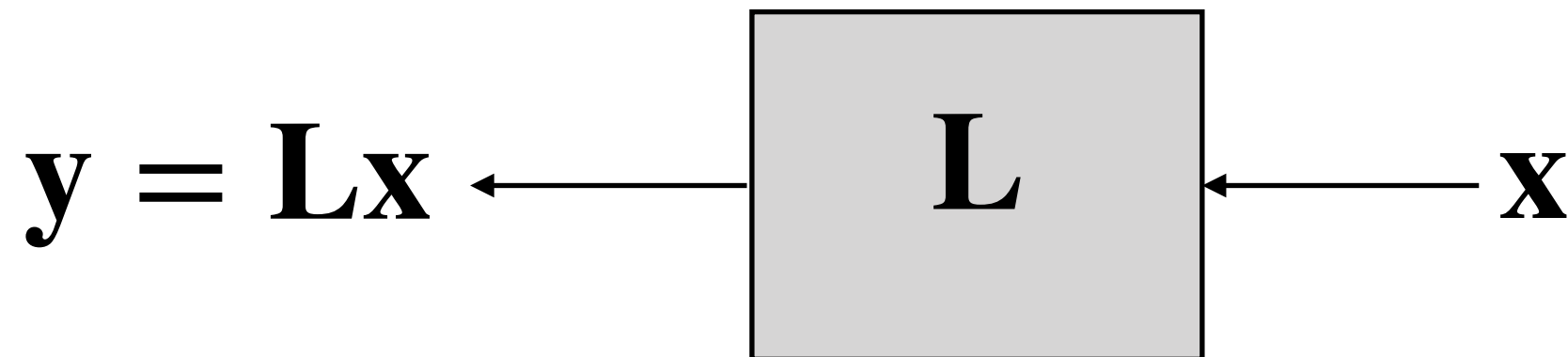
End Do

f: Forward

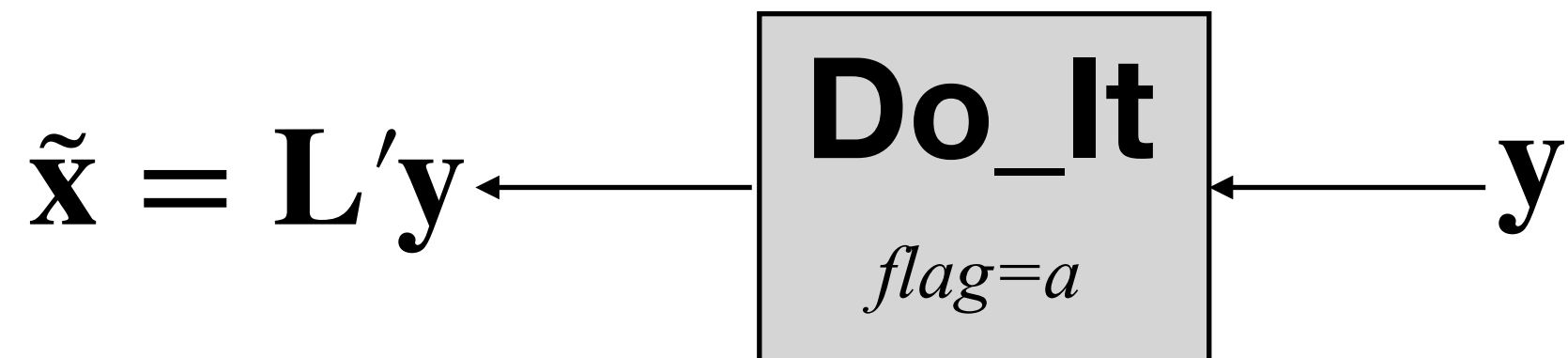
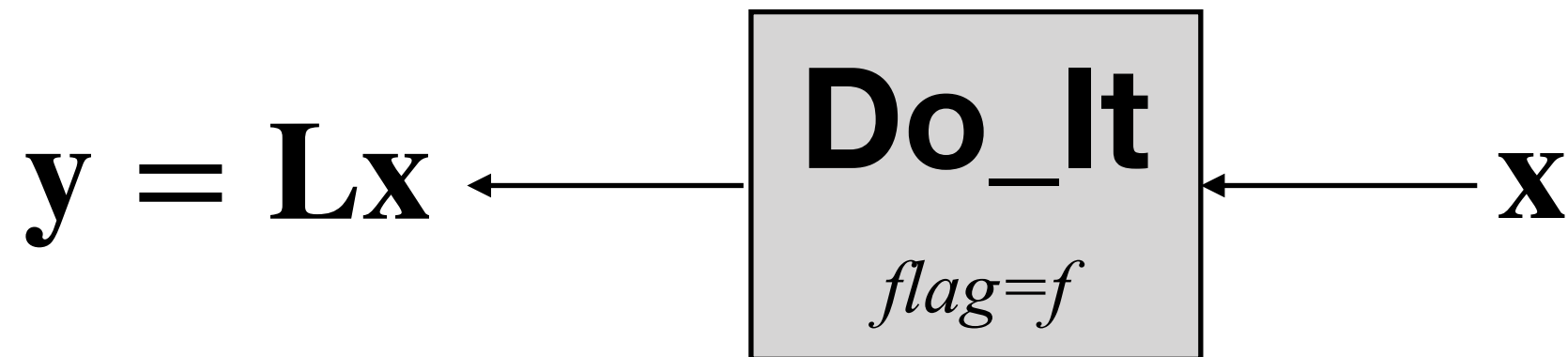
a: Adjoint or transpose

Operators *on the flight*

(Matrices are replaced by implicit form linear operators!!)



Operators *on the flight*



f: Forward

a: Adjoint or transpose

Why this is important?

- Migration and de-migration operators cannot be written as matrices. Least-squares migration requires the **Do_It** approach.
- In addition, time-domain Radon transforms can not be written in explicit form. For Radon problems I also use the **Do_It** approach.
- Reconstruction problems often entail using operators that cannot be written via explicit matrices (FFTs, Curvelet Frames, etc). In this case, we also adopt the **Do_It** approach
- In general, multidimensional problems where \mathbf{m} and \mathbf{d} are not vectors can still be analyzed via linear algebra tools by providing simple rules that make linear operators behave like matrices

Dot-product test

- How do you guarantee that the codes for the forward and adjoint operators behave like \mathbf{L} and \mathbf{L}' ?
- Let's \mathbf{L} be the forward operator and let's call \mathbf{B} the tentative adjoint or transpose operator

$$\mathbf{y}_1 = \mathbf{L} \mathbf{x}_1, \quad \mathbf{x}_2 = \mathbf{B} \mathbf{y}_2$$

- Form the two inner products

$$\mathbf{y}_2^T \mathbf{y}_1 = \mathbf{y}_2 \mathbf{L} \mathbf{x}_1$$

$$\mathbf{x}_1^T \mathbf{x}_2 = \mathbf{x}_1 \mathbf{B} \mathbf{y}_2$$

$$\langle \mathbf{y}_2, \mathbf{y}_1 \rangle = \langle \mathbf{y}_2, \mathbf{L} \mathbf{x}_1 \rangle$$

$$\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \langle \mathbf{x}_1, \mathbf{B} \mathbf{y}_2 \rangle$$

More general notation

Dot-product test

- Notice, the definition of the inner product can be extended to multidimensional array

- For vectors $\mathbf{y}_2^T \mathbf{y}_1 = \mathbf{sum}(y_2. * y_1)$

- For ND-arrays $\langle \mathbf{y}_2, \mathbf{y}_1 \rangle = \mathbf{sum}(y_2[:, :] . * y_1[:, :])$

Dot-product test

The two inner products are equal

$$\mathbf{y}_2^T \mathbf{y}_1 = \mathbf{y}_2 \mathbf{L} \mathbf{x}_1 \qquad \mathbf{x}_1^T \mathbf{x}_2 = \mathbf{x}_1 \mathbf{B} \mathbf{y}_2$$

If and only if $\mathbf{B} = \mathbf{L}'$

Therefore, one can write a code for \mathbf{L} and a code for \mathbf{L}' , do the dot-product test and if the two product are equal then one can say that \mathbf{L} and \mathbf{L}' (packed by the function **Do_It**) behave like a matrix and transpose multiplications, respectively. Then you can safely use all you know about linear algebra to solve an inverse problem!

Dot-product test in practice

```
using LinearAlgebra, FFTW
```

```
# Dot product test example using operators rather than matrices
```

```
# Fourier DFT matrices and its Hermitian Transpose are replaced by
```

```
# on-the-flight FFTs
```

```
M = 512
```

```
x1 = randn(M)
```

```
y1 = fft(x1)
```

```
y2 = randn(M)
```

```
x2 = M*ifft(y2)           => Why I have multiplied by M?
```

```
dot_x = x1'*x2
```

```
dot_y = y1'*y2
```

```
println(dot_x)
```

```
println(dot_y)
```

Conjugate Gradients

- It is more efficient to use the Conjugate Gradient (CG) method than SD.
- I will not discuss CG but it basically amounts to also applying on the flight the operators \mathbf{L} and \mathbf{L}' in each iteration (step)
- The Conjugate Gradient algorithm is quite popular so you can find one in C, Fortran, Matlab, Python, etc by just doing a google search.
- *CG minimizes the general quadratic cost function:*

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2$$

Concatenation of operators

- Forward

$$\mathbf{L} = \mathbf{ABC}$$

- Adjoint

$$\mathbf{L}' = \mathbf{C}'\mathbf{B}'\mathbf{A}'$$

- Dot product test must work for individual operators (codes) and then it will work for \mathbf{L} and \mathbf{L}'

Preconditioning

- When using iterative methods (SD or CG), I prefer not to worry about matrices of weights

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W}\mathbf{m}\|_2^2$$

- Therefore, I like to use the following change of variables

$$\mathbf{W}\mathbf{m} = \mathbf{u} \rightarrow \mathbf{P}\mathbf{u} = \mathbf{m}$$

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{P}\mathbf{u}\|_2^2 + \mu \|\mathbf{u}\|_2^2$$

Preconditioning

- Minimize with iterative solver

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{P}\mathbf{u}\|_2^2 + \mu \|\mathbf{u}\|_2^2$$

- With **Do_It** operators:

Forward : LP

Adjoint : P'L'

Preconditioning

- In this example, I use SD to minimize

$$J = \|\mathbf{d} - \mathbf{LPu}\|_2^2 + \mu \|\mathbf{u}\|_2^2$$

Do until convergence

$$\mathbf{r}^k = (\mathbf{LPu}^k - \mathbf{d})$$

$$\mathbf{u}^{k+1} = \mathbf{u}^k - \alpha(\mathbf{P}'\mathbf{L}'(\mathbf{r}^k) + \mu \mathbf{u}^k)$$

End

$$\mathbf{m}_{sol} = \mathbf{Pu}$$

Preconditioning

- **W** and **P** should behave like the inverse of each other

$$\mathbf{W}\mathbf{m} = \mathbf{u} \rightarrow \mathbf{P}\mathbf{u} = \mathbf{m}$$

Then

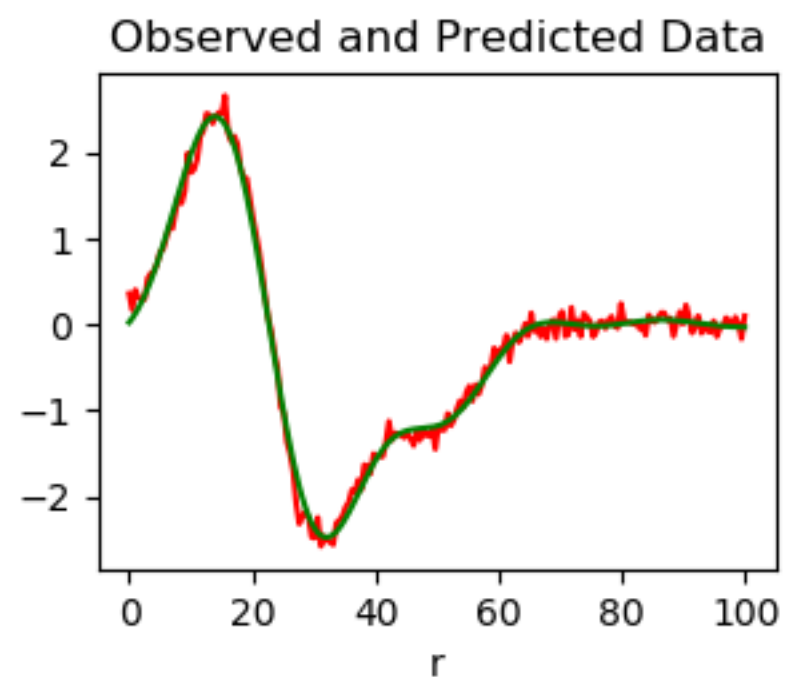
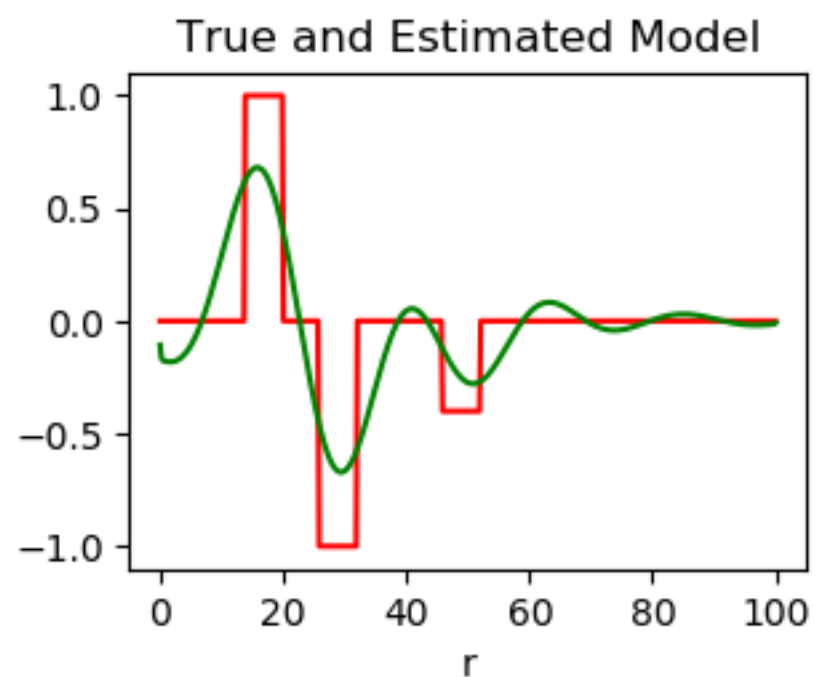
W is a high-pass operator (**applies roughening**)

P is low-pass operator (**applies smoothing**)

Oldie but goodie: S. Ronen, D. Nichols, R. Bale, and R. Ferber, “Dealiasing DMO: Good-pass, bad-pass, and unconstrained”, SEG Technical Program Expanded Abstracts 1995. January 1995, 743-746

Preconditioning (SD)

$$\mathbf{P} = \frac{1}{2} \begin{pmatrix} 2 & 1 & 0 & 0 & 0 \\ 1 & 2 & 1 & 0 & 0 \\ 0 & 1 & 2 & 1 & 0 \\ 0 & 0 & 1 & 2 & 1 \\ 0 & 0 & 0 & 1 & 2 \end{pmatrix}$$



My favourite CG (CGLS)

- I prefer to use a CG algorithm that minimizes the so called Standard Form

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{m}\|_2^2$$

```
function ConjugateGradients(d,operators,parameters;Niter=10,mu=0)
```

```
    r = copy(d)
```

```
    g = LinearOperator(r,operators,parameters,adj=true)
```

$$\mathbf{g} = \mathbf{L}'\mathbf{r}$$

```
    m = zero(g)
```

```
    s = copy(g)
```

```
    gamma = InnerProduct(g,g)
```

```
    for iter = 1 : Niter
```

```
        t = LinearOperator(s,operators,parameters,adj=false)
```

$$\mathbf{t} = \mathbf{L}\mathbf{s}$$

```
        delta = InnerProduct(t,t) + mu*InnerProduct(s,s)
```

```
        alpha = gamma/delta
```

```
        m = m + alpha*s
```

```
        r = r - alpha*t
```

```
        g = LinearOperator(r,operators,parameters,adj=true)
```

$$\mathbf{g} = \mathbf{L}'\mathbf{r}$$

```
        g = g - mu*m
```

```
        gamma0 = copy(gamma)
```

```
        gamma = InnerProduct(g,g)
```

```
        beta = gamma/gamma0
```

```
        s = beta*s + g
```

```
    end
```

```
    return m
```

```
end
```

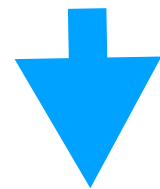
<https://github.com/SeismicJulia/SeisReconstruction.jl/blob/master/src/Tools/ConjugateGradients.jl>

Modified from: <https://web.stanford.edu/group/SOL/software/cgls/>

Advantages of CGLS to minimize the standard form

You can turn a general problem with quadratic regularization into its standard form and use the previous page CGLS algorithm

$$J = \|\mathbf{d} - \mathbf{L}\mathbf{m}\|_2^2 + \mu \|\mathbf{W}\mathbf{m}\|_2^2$$



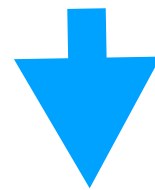
$$J = \|\mathbf{d} - \mathbf{L}\mathbf{P}\mathbf{u}\|_2^2 + \mu \|\mathbf{u}\|_2^2, \text{ with, } \mathbf{u} = \mathbf{P}\mathbf{m}$$

$$\mathbf{P} = \mathbf{W}^{-1} \leftarrow \text{to discuss}$$

Advantages of CGLS to minimize the standard form

You can also concatenate linear operators

$$J = \|\mathbf{W}_d[\mathbf{d} - \mathbf{A}\mathbf{B}\mathbf{m}]\|_2^2 + \mu\|\mathbf{W}_m\mathbf{m}\|_2^2$$



$$J = \|\mathbf{d}_{w_d} - \mathbf{W}_d\mathbf{A}\mathbf{B}\mathbf{P}_m\mathbf{u}\|_2^2 + \mu\|\mathbf{u}\|_2^2, \text{ with, } \mathbf{u} = \mathbf{P}_m\mathbf{m}, \mathbf{d}_{w_d} = \mathbf{W}_d\mathbf{d}$$

$$\text{operators} = [\mathbf{W}_d, \mathbf{A}, \mathbf{B}, \mathbf{P}_m]$$

$$\text{LinearOperator}(\text{Input}, \text{operators}, \text{parameters}, \text{adj}=\text{false}) \quad \text{Output} = \mathbf{W}_d\mathbf{A}\mathbf{B}\mathbf{P}_m\text{Input}$$

$$\text{LinearOperator}(\text{Input}, \text{operators}, \text{parameters}, \text{adj}=\text{true}) \quad \text{Output} = \mathbf{P}_m'\mathbf{B}'\mathbf{A}'\mathbf{W}_d'\text{Input}$$

Example: Time-domain Parabolic Radon Transform

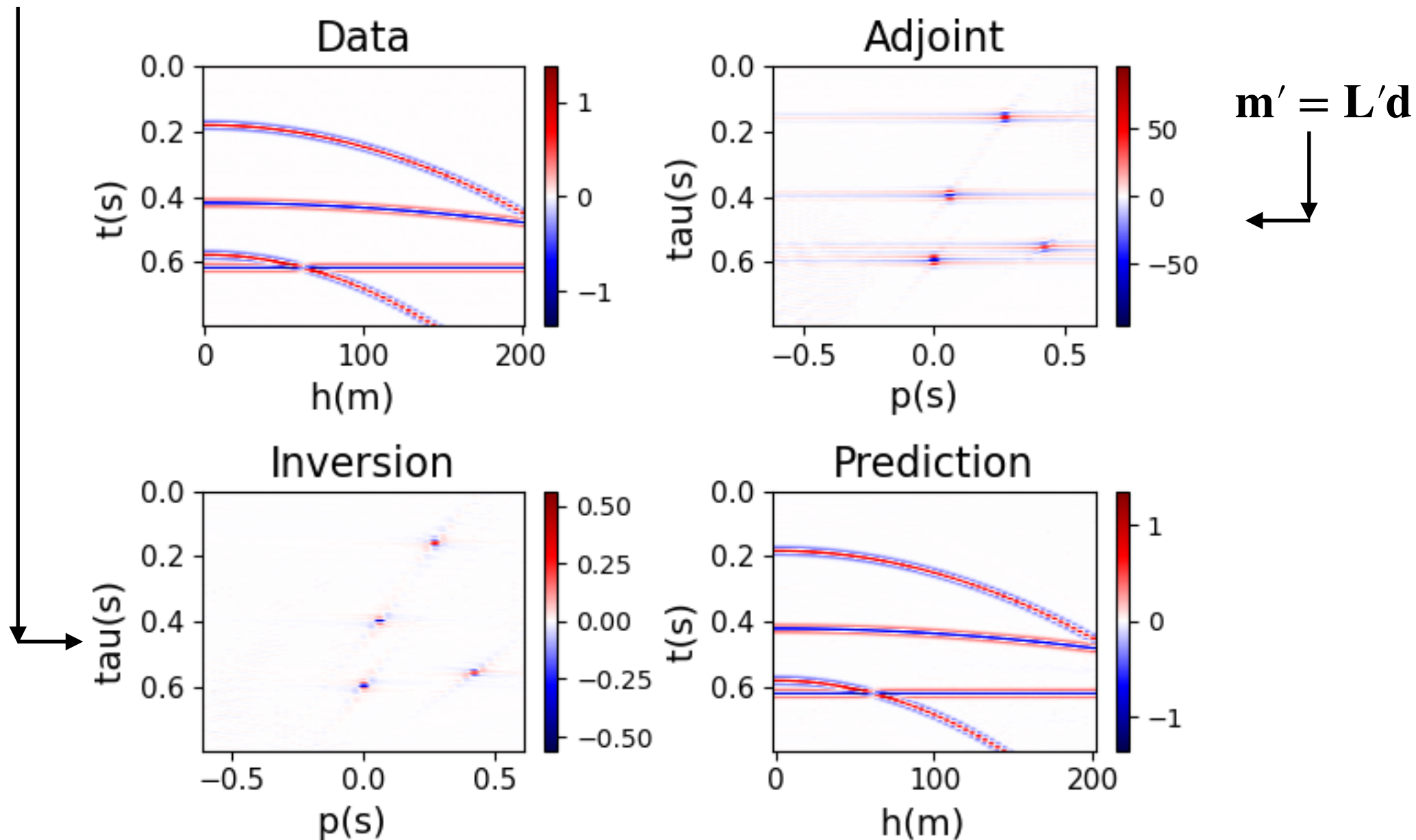
- **Notebook:** `ConjugateGradient_in_SeisReconstruction.ipynbin`
- Compute reflections with parabolic moveout and then use CG to retrieve the Radon domain coefficient that model the reflections.
- Radon adjoint and forward operators:

$$m'(\tau, p) = \sum_h d(t = \tau + pg(h), h) \rightarrow \mathbf{m}' = \mathbf{L}'\mathbf{d}$$
$$d(t, h) = \sum_p m(\tau = t - pg(h)) \rightarrow \mathbf{d} = \mathbf{L}\mathbf{m}$$

- **Linear RT** $g(h) = h, p$: **Slope**
- **Parabolic RT** $g(h) = (h/h_{max})^2, p$: **Residual moveout at far offset**

Example: Time-domain Parabolic Radon Transform

$$\mathbf{m}_{sol} = \operatorname{argmin}_{\mathbf{m}} \|\mathbf{L}\mathbf{m} - \mathbf{d}\|_2^2 + \mu \|\mathbf{m}\|_2^2 \quad \text{Via CGLS}$$



CGLS_versus_SD_example.ipynb

Comparison of SD and CGLS path for a simple problem with two unknown where one can visualize the cost function

