

Gauss–Newton and Levenberg–Marquardt Algorithms

Roland Hostettler

September 26, 2018

Recap

- ▶ Sequential linear least squares:

$$\hat{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_{n-1} + \mathbf{L}_n(\mathbf{y}_n - \mathbf{C}_n\hat{\boldsymbol{\theta}}_{n-1})$$

$$\mathbf{P}_n = \text{Cov}\{\hat{\boldsymbol{\theta}}_n\} = \mathbf{P}_{n-1} - \mathbf{L}_n(\mathbf{C}_n\mathbf{P}_{n-1}\mathbf{C}_n^\top + \mathbf{R}_n)\mathbf{L}_n^\top$$

$$\mathbf{L}_n = \mathbf{P}_{n-1}\mathbf{C}_n^\top(\mathbf{C}_n\mathbf{P}_{n-1}\mathbf{C}_n^\top + \mathbf{R}_n)^{-1}$$

- ▶ Regularized linear least squares:

$$\hat{\boldsymbol{\theta}} = \mathbf{m} + \mathbf{K}(\mathbf{y} - \mathbf{G}\mathbf{m})$$

$$\text{Cov}\{\hat{\boldsymbol{\theta}}\} = \mathbf{P} - \mathbf{K}(\mathbf{G}\mathbf{P}\mathbf{G}^\top + \mathbf{R})\mathbf{K}^\top$$

$$\mathbf{K} = \mathbf{P}\mathbf{G}^\top(\mathbf{G}\mathbf{P}\mathbf{G}^\top + \mathbf{R})^{-1}$$

- ▶ Gradient descent:

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} - \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}^{(i)}}$$

$$\Delta\boldsymbol{\theta}^{(i+1)} = -\mathbf{G}_{\boldsymbol{\theta}}^\top \mathbf{R}^{-1}(\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))$$

Intended Learning Outcomes

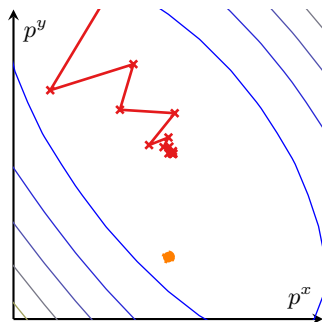
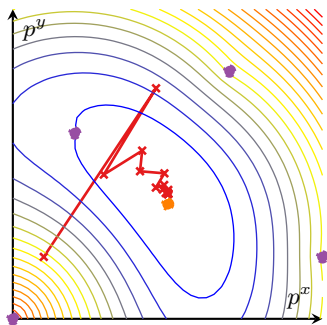
After this lecture, you will be able to:

- ▶ Describe the Gauss–Newton and Levenberg–Marquardt algorithms;
- ▶ discuss the properties of the gradient descent, Gauss–Newton, and Levenberg–Marquardt algorithms;
- ▶ recall suitable ways of evaluating the convergence of numerical optimization methods.

Example: Localizing a Target

- ▶ Parameters: Target position $\theta = [p_t^x \ p_t^y]^\top$
- ▶ Measurements model (with $\Delta \mathbf{p}_n = \mathbf{p}_t - \mathbf{p}_n$):

$$y_n = |\Delta \mathbf{p}_n| + r_n = \sqrt{(\Delta p_n^x)^2 + (\Delta p_n^y)^2} + r_n$$



Gauss–Newton Algorithm: Derivation (1/3)

- ▶ Idea: Given $\hat{\boldsymbol{\theta}}^{(i)}$, we can linearize the nonlinear measurement model around that point
- ▶ Linearized measurement model:

$$g(\boldsymbol{\theta}) \approx g(\hat{\boldsymbol{\theta}}^{(i)}) + \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})$$

- ▶ Cost function approximation:

$$\begin{aligned} J_{\text{WLS}}(\boldsymbol{\theta}) &= (\mathbf{y} - g(\boldsymbol{\theta}))^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\boldsymbol{\theta})) \\ &\approx (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}) - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}))^{\top} \mathbf{R}^{-1} \\ &\quad \times (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}) - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})) \\ &= (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}))^{\top} \mathbf{R}^{-1} (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})) \end{aligned}$$

Gauss–Newton Algorithm: Derivation (2/3)

- ▶ Cost function approximation:

$$J_{\text{WLS}}(\boldsymbol{\theta}) \approx (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}))^{\top} \mathbf{R}^{-1} (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}))$$

- ▶ Gradient of the cost function approximation w.r.t. $\boldsymbol{\theta}$:

$$\begin{aligned} \frac{\partial J_{\text{WLS}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &\approx \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}))^{\top} \mathbf{R}^{-1} (\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})) \\ &= -2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)} + 2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \end{aligned}$$

- ▶ Solving for $\boldsymbol{\theta}$:

$$\begin{aligned} \boldsymbol{\theta} &= \hat{\boldsymbol{\theta}}^{(i)} + (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)} \\ &= \hat{\boldsymbol{\theta}}^{(i)} + (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)})) \end{aligned}$$

Gauss–Newton Algorithm: Derivation (3/3)

- ▶ Solution of the **linearized** cost function:

$$\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}^{(i)} + (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))$$

- ▶ Gauss–Newton update:

$$\begin{aligned}\hat{\boldsymbol{\theta}}^{(i+1)} &= \hat{\boldsymbol{\theta}}^{(i)} + \gamma \Delta \hat{\boldsymbol{\theta}}^{(i+1)}, \\ \Delta \hat{\boldsymbol{\theta}}^{(i+1)} &= (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))\end{aligned}$$

- ▶ The step length γ can be selected using a line search with step direction $\Delta \hat{\boldsymbol{\theta}}^{(i+1)}$

Gauss–Newton Algorithm: Algorithm

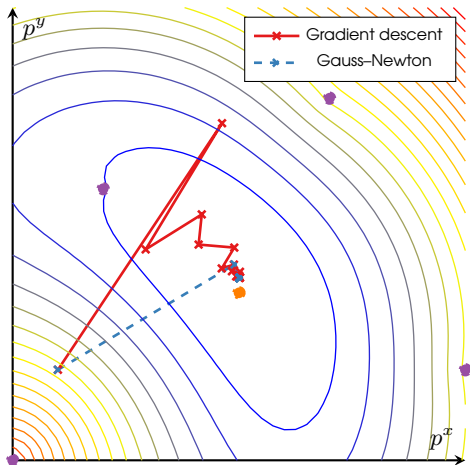
Algorithm 1 Gauss–Newton Algorithm with Line Search

Input: Initial guess $\hat{\theta}^{(0)}$, data \mathbf{y} , function $g(\theta)$, Jacobian \mathbf{G}_θ

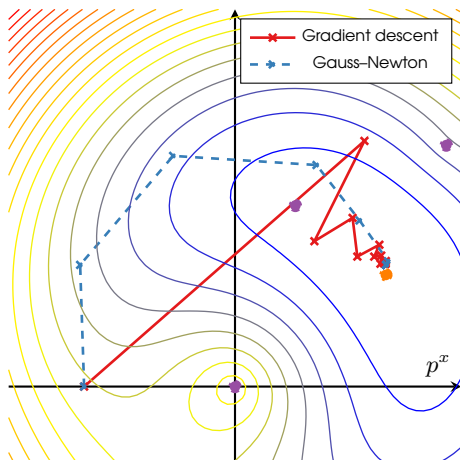
Output: Parameter estimate $\hat{\theta}_{\text{WLS}}$

- 1: Set $i \leftarrow 0$
 - 2: **repeat**
 - 3: Calculate $\Delta\theta^{(i+1)} = (\mathbf{G}_\theta^T \mathbf{R}^{-1} \mathbf{G}_\theta)^{-1} \mathbf{G}_\theta^T \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\theta}^{(i)}))$
 - 4: Set $\gamma^{(i+1)} \leftarrow 1$
 - 5: **repeat** ▷ Line Search
 - 6: Calculate $\hat{\theta}^{(i+1)} = \hat{\theta}^{(i)} + \gamma^{(i+1)} \Delta\theta^{(i+1)}$
 - 7: Set $\gamma^{(i+1)} \leftarrow \gamma^{(i+1)} / 2$
 - 8: **until** $J_{\text{WLS}}(\hat{\theta}^{(i+1)}) < J_{\text{WLS}}(\hat{\theta}^{(i)})$
 - 9: Set $i \leftarrow i + 1$
 - 10: **until** Converged
-

Example: Localizing a Target (2)



Example: Localizing a Target (3)



Levenberg–Marquardt Algorithm: Motivation

- ▶ Gradient descent: Quickly moves to low cost area, creeps to minimum
- ▶ Gauss–Newton: Straight to minimum, may take a detour
- ▶ Can we have the best of both worlds?
- ▶ ... kind of, the Levenberg–Marquardt algorithm

Levenberg–Marquardt Algorithm: Derivation (1/2)

- ▶ The Levenberg–Marquardt algorithm can be seen as a regularized version of the Gauss–Newton algorithm
- ▶ Cost function approximation:

$$J(\boldsymbol{\theta}) \approx \left(\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \right)^{\top} \mathbf{R}^{-1} \left(\mathbf{e}^{(i)} - \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \right)^{\top} + \lambda(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})^{\top}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})$$

with $\mathbf{e}^{(i)} = \mathbf{y} - g(\boldsymbol{\theta}^{(i)})$

- ▶ Gradient of the cost function:

$$\begin{aligned} \frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &\approx -2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)} + 2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) + 2\lambda(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \\ &= -2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)} + 2(\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}} + \lambda \mathbf{I})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \end{aligned}$$

Levenberg–Marquardt Algorithm: Derivation (2/2)

- ▶ Gradient of the cost function:

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \approx -2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)} + 2(\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}} + \lambda \mathbf{I})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})$$

- ▶ Solving for $\boldsymbol{\theta}$:

$$\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}^{(i)} + (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}} + \lambda \mathbf{I})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{e}^{(i)}$$

- ▶ Levenberg–Marquardt update:

$$\begin{aligned} \hat{\boldsymbol{\theta}}^{(i+1)} &= \hat{\boldsymbol{\theta}}^{(i)} + \Delta \hat{\boldsymbol{\theta}}^{(i+1)}, \\ \Delta \hat{\boldsymbol{\theta}}^{(i+1)} &= (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}} + \lambda \mathbf{I})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)})). \end{aligned}$$

Levenberg–Marquardt Algorithm: Damping (1/2)

- ▶ Levenberg–Marquardt update:

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} + \Delta \hat{\boldsymbol{\theta}}^{(i+1)},$$
$$\Delta \hat{\boldsymbol{\theta}}^{(i+1)} = (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}} + \lambda \mathbf{I})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)})).$$

- ▶ How should the *damping parameter* λ be chosen?
- ▶ If $\lambda \rightarrow 0$:

$$\Delta \hat{\boldsymbol{\theta}}^{(i+1)} \rightarrow (\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} \mathbf{G}_{\boldsymbol{\theta}})^{-1} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))$$

- ▶ If $\lambda \rightarrow \infty$:

$$\Delta \hat{\boldsymbol{\theta}}^{(i+1)} \rightarrow \frac{1}{\lambda} \mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)})).$$

Levenberg–Marquardt Algorithm: Damping (2/2)

- ▶ The predicted (linearized) reduction in cost is:

$$\begin{aligned}\Delta J(\boldsymbol{\theta}) &= J(\hat{\boldsymbol{\theta}}^{(i)}) - J(\boldsymbol{\theta}) \\ &\approx 2(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})^\top \mathbf{G}_\theta^\top \mathbf{R}^{-1}(\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)})) \\ &\quad - (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)})^\top \mathbf{G}_\theta^\top \mathbf{R}^{-1} \mathbf{G}_\theta (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{(i)}) \\ &\triangleq \Delta \tilde{J}(\boldsymbol{\theta})\end{aligned}$$

- ▶ The relative gain in cost is then:

$$\rho(\boldsymbol{\theta}) = \frac{J(\hat{\boldsymbol{\theta}}^{(i)}) - J(\boldsymbol{\theta})}{\Delta \tilde{J}(\boldsymbol{\theta})}$$

- ▶ Adaption strategy:

- ▶ If $\rho > 0$, accept $\Delta \hat{\boldsymbol{\theta}}^{(i+1)}$ and decrease the damping

$$\lambda \leftarrow \lambda \max(1/3, 1 - (2\rho - 1)^3), \quad \nu = 2,$$

- ▶ otherwise, discard $\Delta \hat{\boldsymbol{\theta}}^{(i+1)}$ increase the damping by setting

$$\lambda \leftarrow \nu \lambda, \quad \nu \leftarrow 2\nu.$$

Levenberg–Marquardt Algorithm: Algorithm

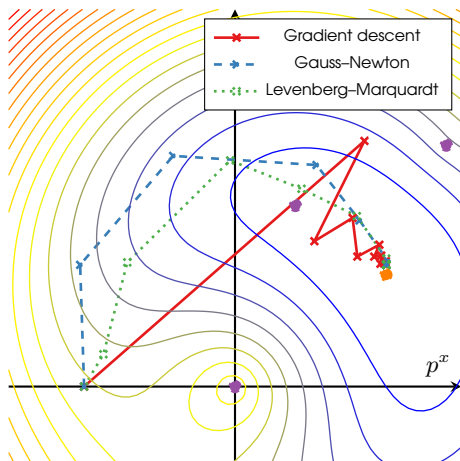
Algorithm 2 Levenberg–Marquardt Algorithm

Input: Initial guess $\hat{\theta}^{(0)}$, data \mathbf{y} , function $g(\theta)$, Jacobian \mathbf{G}_θ

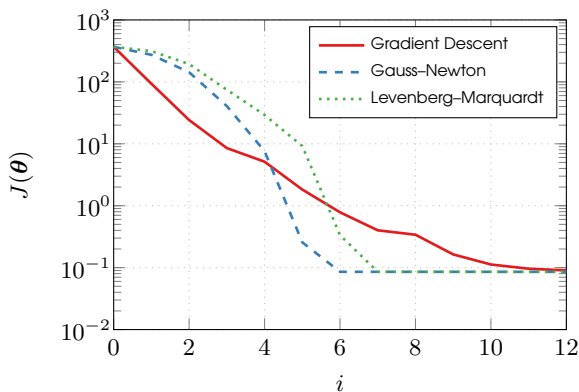
Output: Parameter estimate $\hat{\theta}_{\text{WLS}}$

- 1: Set $i \leftarrow 0$, $\lambda = \max(\text{diag}(\mathbf{G}_\theta^T \mathbf{R}^{-1} \mathbf{G}_\theta))$, and $\nu \leftarrow 2$
- 2: **repeat**
- 3: **repeat**
- 4: Calculate $\Delta \hat{\theta}^{(i+1)}$ and ρ
- 5: **if** $\rho > 0$ **then**
- 6: Set $\hat{\theta}^{(i+1)} = \hat{\theta}^{(i+1)} + \Delta \hat{\theta}^{(i+1)}$
- 7: Set $\lambda \leftarrow \lambda \max(1/3, 1 - (2\rho - 1)^3)$, $\nu = 2$
- 8: **else**
- 9: Set $\lambda \leftarrow \nu \lambda$ and $\nu \leftarrow 2\nu$
- 10: **end if**
- 11: **until** $\Delta \theta^{(i+1)}$ is accepted
- 12: Set $i \leftarrow i + 1$
- 13: **until** Converged

Example: Localizing a Target (4)



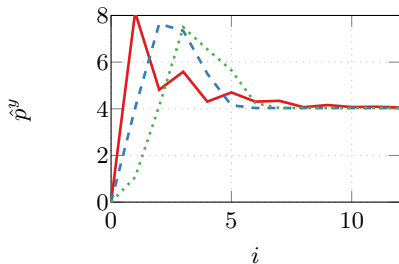
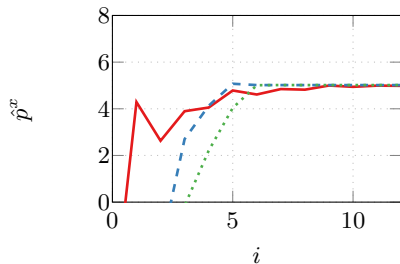
Convergence Criteria: Decrease in Cost



- Relative change in cost:

$$\frac{J(\hat{\theta}^{(i+1)}) - J(\hat{\theta}^{(i)})}{J(\hat{\theta}^{(i)})}$$

Convergence Criteria: Change in Parameters



- Relative change in parameter magnitude:

$$\frac{|\hat{\theta}^{(i+1)} - \hat{\theta}^{(i)}|}{|\hat{\theta}^{(i)}|}$$

Convergence Criteria: Overview

Terminate the numerical search when:

- ▶ The relative change in cost

$$\frac{J(\hat{\theta}^{(i+1)}) - J(\hat{\theta}^{(i)})}{J(\hat{\theta}^{(i)})}$$

falls below a threshold ϵ ,

- ▶ the relative change in parameter magnitude

$$\frac{|\hat{\theta}^{(i+1)} - \hat{\theta}^{(i)}|}{|\hat{\theta}^{(i)}|}$$

falls below a threshold ϵ ,

- ▶ the maximum number of iterations I_{\max} is reached.

Summary

- ▶ Gauss–Newton algorithm:

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} + \gamma(\mathbf{G}_{\hat{\boldsymbol{\theta}}^i}^T \mathbf{R}^{-1} \mathbf{G}_{\hat{\boldsymbol{\theta}}^i})^{-1} \mathbf{G}_{\hat{\boldsymbol{\theta}}^i}^T \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))$$

- ▶ Levenberg–Marquardt algorithm:

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} + (\mathbf{G}_{\hat{\boldsymbol{\theta}}^i}^T \mathbf{R}^{-1} \mathbf{G}_{\hat{\boldsymbol{\theta}}^i} + \lambda \mathbf{I})^{-1} \mathbf{G}_{\hat{\boldsymbol{\theta}}^i}^T \mathbf{R}^{-1} (\mathbf{y} - g(\hat{\boldsymbol{\theta}}^{(i)}))$$

