

Sequential and Regularized Least Squares, Static Nonlinear Models and Gradient Descent

Roland Hostettler

September 13, 2018

Recap

- ▶ The general linear model is given by

$$\mathbf{y} = \mathbf{G}\boldsymbol{\theta} + \mathbf{r}, \quad \mathbb{E}\{\mathbf{r}\} = 0, \quad \text{Cov}\{\mathbf{r}\} = \mathbf{R}$$

- ▶ The weighted linear least squares estimator is

$$\hat{\boldsymbol{\theta}}_{WLS} = (\mathbf{G}^T \mathbf{R}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \mathbf{R}^{-1} \mathbf{y}$$

with $\mathbb{E}\{\hat{\boldsymbol{\theta}}_{WLS}\} = \boldsymbol{\theta}$ and $\text{Cov}\{\hat{\boldsymbol{\theta}}_{WLS}\} = (\mathbf{G}^T \mathbf{R}^{-1} \mathbf{G})^{-1}$.

- ▶ The LS estimator is a special case of the WLS squares estimator with $\mathbf{R} \propto \mathbf{I}$
- ▶ It holds that

$$\text{Cov}\{\hat{\boldsymbol{\theta}}_{WLS}\} \leq \text{Cov}\{\hat{\boldsymbol{\theta}}_{LS}\}$$

Intended Learning Outcomes

After this lecture, you will be able to:

- ▶ Describe and identify sequential and regularized linear least squares estimators;
- ▶ distinguish (weighted) linear least squares and regularized linear least squares estimators and their properties;
- ▶ identify the challenges encountered in nonlinear sensor models;
- ▶ restate the gradient descent algorithm for nonlinear least squares.

Sequential Linear Least Squares: Formulation

- ▶ Sometimes data arrives sequentially; how can we update an existing estimate with the new data?
- ▶ Given the data $\mathbf{y}_{1:n-1} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n-1}\}$, we have

$$\hat{\boldsymbol{\theta}}_{n-1} = (\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{R}^{-1} \mathbf{y}$$
$$\text{Cov}\{\hat{\boldsymbol{\theta}}_{n-1}\} = (\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{G})^{-1} \triangleq \mathbf{P}_{n-1}$$

- ▶ Updated cost function when \mathbf{y}_n arrives:

$$J_{\text{SLS}}(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{G}\boldsymbol{\theta})^\top \mathbf{R}^{-1} (\mathbf{y} - \mathbf{G}\boldsymbol{\theta})$$
$$+ (\mathbf{y}_n - \mathbf{C}_n \boldsymbol{\theta})^\top \mathbf{R}_n^{-1} (\mathbf{y}_n - \mathbf{C}_n \boldsymbol{\theta}),$$

Sequential Linear Least Squares: Derivation (1/2)

- ▶ Updated cost function when \mathbf{y}_n arrives:

$$J_{\text{SLS}}(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{G}\boldsymbol{\theta})^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{G}\boldsymbol{\theta}) \\ + (\mathbf{y}_n - \mathbf{C}_n\boldsymbol{\theta})^T \mathbf{R}_n^{-1}(\mathbf{y}_n - \mathbf{C}_n\boldsymbol{\theta}),$$

- ▶ Gradient:

$$\frac{\partial J_{\text{SLS}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2(\mathbf{G}^T \mathbf{R}^{-1} \mathbf{y} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{y}_n) \\ + 2(\mathbf{G}^T \mathbf{R}^{-1} \mathbf{G} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{C}_n) \boldsymbol{\theta}$$

- ▶ Updated least squares estimator:

$$\hat{\boldsymbol{\theta}}_n = (\mathbf{P}_{n-1}^{-1} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{C}_n)^{-1} (\mathbf{G}^T \mathbf{R}^{-1} \mathbf{y} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{y}_n).$$

Sequential Linear Least Squares: Derivation (2/2)

- ▶ Updated least squares estimator:

$$\hat{\boldsymbol{\theta}}_n = (\mathbf{P}_{n-1}^{-1} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{C}_n)^{-1} (\mathbf{G}^T \mathbf{R}^{-1} \mathbf{y} + \mathbf{C}_n^T \mathbf{R}_n^{-1} \mathbf{y}_n).$$

⇒ Problem: Does not re-use $\hat{\boldsymbol{\theta}}_{n-1}$ and \mathbf{P}_{n-1} well

- ▶ Alternative form:

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

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

Sequential Linear Least Squares: Properties

- ▶ Sequential linear least squares:

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

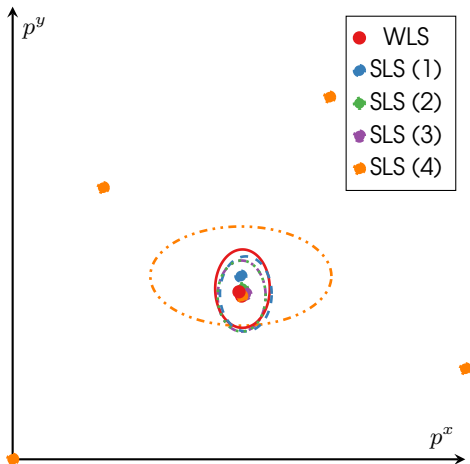
- ▶ Mean:

$$\begin{aligned} \mathbb{E}\{\hat{\boldsymbol{\theta}}_n\} &= \mathbb{E}\{\hat{\boldsymbol{\theta}}_{n-1}\} + \mathbf{L}_n(\mathbb{E}\{\mathbf{y}_n\} - \mathbf{C}_n \mathbb{E}\{\hat{\boldsymbol{\theta}}_{n-1}\}) \\ &= \boldsymbol{\theta} + \mathbf{L}_n(\mathbf{C}_n \boldsymbol{\theta} - \mathbf{C}_n \boldsymbol{\theta}) \\ &= \boldsymbol{\theta}, \end{aligned}$$

- ▶ Covariance:

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

Example: Localizing a Target (1)



Regularized Linear Least Squares: Formulation

- ▶ Sometimes we have prior knowledge about θ
- ▶ For example, we could assume that

$$E\{\theta\} = m \text{ and } \text{Cov}\{\theta\} = P$$

- ▶ Cost function with regularization term:

$$J_{\text{ReLS}}(\theta) = \underbrace{(y - G\theta)^T R^{-1} (y - G\theta)}_{\text{Data}} + \underbrace{(\theta - m)^T P^{-1} (\theta - m)}_{\text{Prior Knowledge}}$$

- ▶ This corresponds to a *Bayesian linear model*

Regularized Linear Least Squares: Derivation

- ▶ Gradient of the cost function:

$$\begin{aligned}\frac{\partial J_{\text{ReLS}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{y} - \mathbf{G}\boldsymbol{\theta})^\top \mathbf{R}^{-1} (\mathbf{y} - \mathbf{G}\boldsymbol{\theta}) \\ &\quad + \frac{\partial}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta} - \mathbf{m})^\top \mathbf{P}^{-1} (\boldsymbol{\theta} - \mathbf{m}) \\ &= -2\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{y} + 2\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{G}\boldsymbol{\theta} - 2\mathbf{P}^{-1} \mathbf{m} + 2\mathbf{P}^{-1} \boldsymbol{\theta} \\ &= -2(\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{y} + \mathbf{P}^{-1} \mathbf{m}) + 2(\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{G} + \mathbf{P}^{-1}) \boldsymbol{\theta}\end{aligned}$$

- ▶ Setting to zero and solving for $\boldsymbol{\theta}$:

$$\hat{\boldsymbol{\theta}}_{\text{ReLS}} = (\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{G} + \mathbf{P}^{-1})^{-1} (\mathbf{G}^\top \mathbf{R}^{-1} \mathbf{y} + \mathbf{P}^{-1} \mathbf{m})$$

- ▶ Alternative formulation (using matrix inversion lemma):

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{\text{ReLS}} &= \mathbf{m} + \mathbf{K}(\mathbf{y} - \mathbf{G}\mathbf{m}) \\ \mathbf{K} &= \mathbf{P}\mathbf{G}^\top (\mathbf{G}\mathbf{P}\mathbf{G}^\top + \mathbf{R})^{-1}\end{aligned}$$

Regularized Linear Least Squares: Properties

- ▶ Regularized linear least squares:

$$\hat{\boldsymbol{\theta}}_{\text{ReLS}} = \mathbf{m} + \mathbf{K}(\mathbf{y} - \mathbf{G}\mathbf{m})$$
$$\mathbf{K} = \mathbf{P}\mathbf{G}^{\top}(\mathbf{G}\mathbf{P}\mathbf{G}^{\top} + \mathbf{R})^{-1}$$

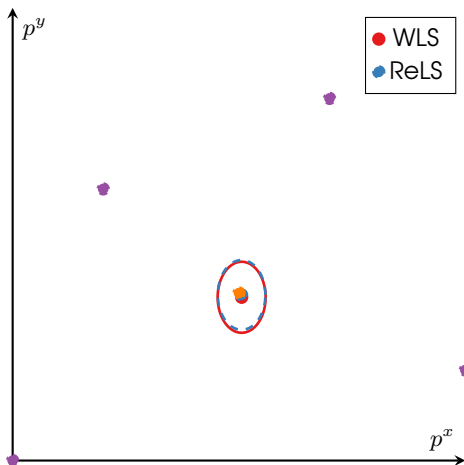
- ▶ Mean:

$$\begin{aligned} \mathbb{E}\{\hat{\boldsymbol{\theta}}_{\text{ReLS}}\} &= \mathbb{E}\{\mathbf{m} + \mathbf{K}(\mathbf{y} - \mathbf{G}\mathbf{m})\} \\ &= \mathbf{m} + \mathbf{K}(\mathbf{G}\mathbb{E}\{\boldsymbol{\theta}\} - \mathbf{G}\mathbf{m}). \end{aligned}$$

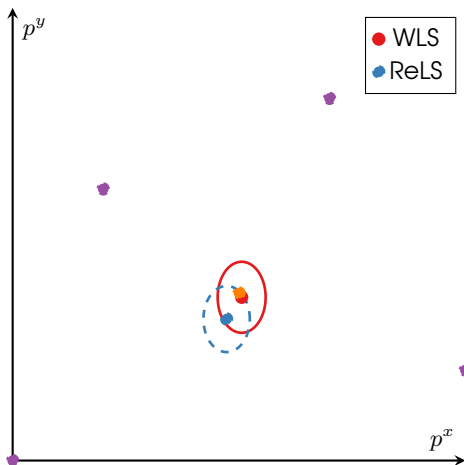
- ▶ Covariance:

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

Example: Localizing a Target (2)



Example: Localizing a Target (3)



Static Nonlinear Models

- ▶ Linear models have closed form solutions, but are limited in many cases
- ▶ General nonlinear model:

$$\mathbf{y} = g(\boldsymbol{\theta}) + \mathbf{r},$$

- ▶ General cost function:

$$J_{\text{WLS}}(\boldsymbol{\theta}) = (\mathbf{y} - g(\boldsymbol{\theta}))^{\text{T}} \mathbf{R}^{-1} (\mathbf{y} - g(\boldsymbol{\theta})).$$

- ▶ For some models, closed form solutions do exist, ...
- ▶ ...but for most they do not

Numerical Optimization

- ▶ Iterative algorithms to find the stationary points or roots of a function
- ▶ Generally find **local minima** \Rightarrow Require good initialization
- ▶ Two predominant approaches:
 1. Methods that find the minima of a function, that is, methods that solve

$$x = \underset{x}{\operatorname{argmin}} f(x)$$

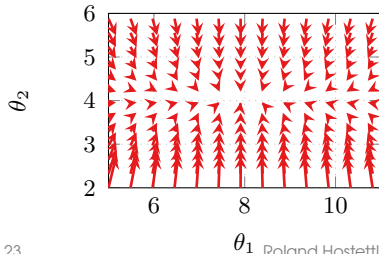
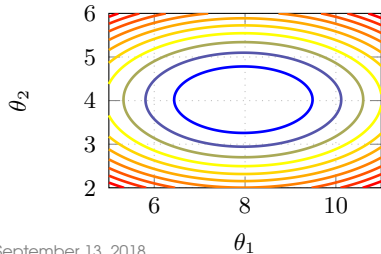
2. Methods that find the zeros of the gradient, that is, methods that solve

$$\frac{\partial f(x)}{\partial x} = 0$$

Gradient Descent: Formulation

- ▶ The *gradient* of $f(x)$ w.r.t. x points to the direction where $f(x)$ increases as a function of x
- ▶ Changing x in the **opposite direction** of the gradient decreases $f(x)$
- ▶ If the function to minimize is $J_{\text{WLS}}(\boldsymbol{\theta})$, the cost is decreased by the iteration

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} - \gamma \nabla_{\boldsymbol{\theta}} J_{\text{WLS}}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}^{(i)}},$$



Gradient Descent: Derivation (1/3)

- ▶ Scalar cost function

$$J_{\text{LS}}(\boldsymbol{\theta}) = \sum_{n=1}^N (y_n - g_n(\boldsymbol{\theta}))^2$$

- ▶ Gradient

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J_{\text{LS}}(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \sum_{n=1}^N (y_n - g_n(\boldsymbol{\theta}))^2 \\ &= \sum_{n=1}^N -2 \nabla_{\boldsymbol{\theta}} g_n(\boldsymbol{\theta}) (y_n - g_n(\boldsymbol{\theta}))\end{aligned}$$

Gradient Descent: Derivation (2/3)

- ▶ Gradient

$$\nabla_{\boldsymbol{\theta}} J_{\text{LS}}(\boldsymbol{\theta}) = \sum_{n=1}^N -2 \nabla_{\boldsymbol{\theta}} g_n(\boldsymbol{\theta}) (y_n - g_n(\boldsymbol{\theta}))$$

- ▶ Vector form:

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J_{\text{LS}}(\boldsymbol{\theta}) &= -2 \begin{bmatrix} \frac{\partial g_1(\boldsymbol{\theta})}{\partial \theta_1} & \frac{\partial g_2(\boldsymbol{\theta})}{\partial \theta_1} & \cdots & \frac{\partial g_N(\boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial g_1(\boldsymbol{\theta})}{\partial \theta_2} & \frac{\partial g_2(\boldsymbol{\theta})}{\partial \theta_2} & & \vdots \\ \vdots & & \ddots & \frac{\partial g_N(\boldsymbol{\theta})}{\partial \theta_{K-1}} \\ \frac{\partial g_1(\boldsymbol{\theta})}{\partial \theta_K} & \cdots & \frac{\partial g_{N-1}(\boldsymbol{\theta})}{\partial \theta_K} & \frac{\partial g_N(\boldsymbol{\theta})}{\partial \theta_K} \end{bmatrix} (\mathbf{y} - g(\boldsymbol{\theta})) \\ &= \mathbf{G}_{\boldsymbol{\theta}}^{\text{T}} (\mathbf{y} - g(\boldsymbol{\theta})) \end{aligned}$$

- ▶ $\mathbf{G}_{\boldsymbol{\theta}}$ is the **Jacobian** matrix of $g(\boldsymbol{\theta})$

Gradient Descent: Derivation (3/3)

- ▶ Generalization to WLS cost function:

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J_{\text{WLS}}(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} (\mathbf{y} - g(\boldsymbol{\theta}))^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\boldsymbol{\theta})) \\ &= -2\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\boldsymbol{\theta}))\end{aligned}$$

- ▶ The **direction of the negative gradient** is

$$-\nabla_{\boldsymbol{\theta}} J_{\text{WLS}}(\boldsymbol{\theta}) = -\mathbf{G}_{\boldsymbol{\theta}}^{\top} \mathbf{R}^{-1} (\mathbf{y} - g(\boldsymbol{\theta}))$$

- ▶ The parameter update becomes:

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

- ▶ The Jacobian matrix $\mathbf{G}_{\boldsymbol{\theta}}$ is evaluated at $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}^{(i)}$

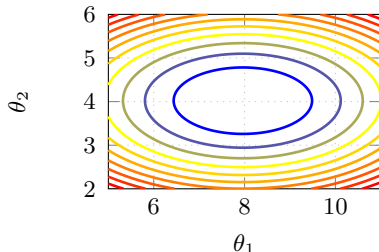
Gradient Descent: Line Search (1/2)

- ▶ Parameter update:

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} + \gamma \Delta \boldsymbol{\theta}^{(i+1)}$$

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

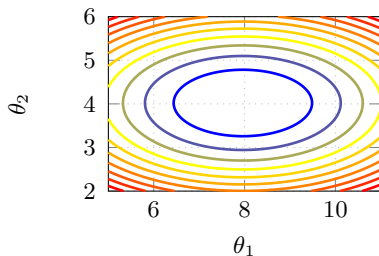
- ▶ How long should γ be?
 - ▶ Large γ : Might overshoot
 - ▶ Small γ : Too slow progress



Gradient Descent: Line Search (2/2)

Adaptive step length selection strategy:

1. Choose $\gamma = 1$
2. Calculate $\hat{\theta}^{(i+1)}$ and the corresponding cost
3. If the cost is lowered, accept $\hat{\theta}^{(i+1)}$, otherwise set $\gamma \leftarrow \gamma/2$ and return to step 2



Gradient Descent: Algorithm

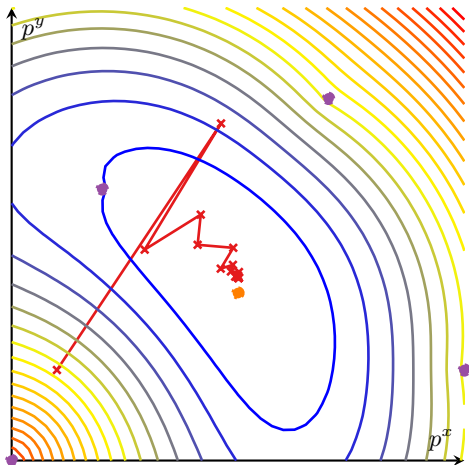
Algorithm 1 Gradient Descent with Line Search

Require: Initial guess $\hat{\theta}^{(0)}$, data \mathbf{y} , function $g(\theta)$, Jacobian G_{θ}

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

- 1: Set $i \leftarrow 0$
 - 2: **repeat**
 - 3: Calculate $\Delta\theta^{(i+1)} = G_{\theta}^T R^{-1}(\mathbf{y} - g(\hat{\theta}^{(i)}))$
 - 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 (4)



Summary

- ▶ The sequential linear least squares estimator is

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

with $\mathbb{E}\{\hat{\boldsymbol{\theta}}_n\} = \boldsymbol{\theta}$ and

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

- ▶ The regularized linear least squares estimator is

$$\hat{\boldsymbol{\theta}}_{\text{ReLS}} = \mathbf{m} + \mathbf{K}(\mathbf{y} - \mathbf{G}\mathbf{m})$$
$$\mathbf{K} = \mathbf{P}\mathbf{G}^\top (\mathbf{G}\mathbf{P}\mathbf{G}^\top + \mathbf{R})^{-1}$$

with $\mathbb{E}\{\hat{\boldsymbol{\theta}}_{\text{ReLS}}\} = \mathbf{m} + \mathbf{K}(\mathbf{G}\mathbb{E}\{\boldsymbol{\theta}\} - \mathbf{G}\mathbf{m})$ and

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

- ▶ The gradient descent algorithm is a numerical method for solving nonlinear least squares problems