

# Least squares regularization

Lecture 08

September 23, 2021

Beck, sections 3.3-3.4

# Least squares problems

L08-S01

Recall:  $\min_{\underline{x} \in \mathbb{R}^N} f(\underline{x}),$

$$f(\underline{x}) = \underline{x}^T \underline{A}^T \underline{A} \underline{x} - 2 \underline{b}^T \underline{A} \underline{x} + \|\underline{b}\|_2^2$$
$$= \|\underline{A} \underline{x} - \underline{b}\|_2^2$$

If  $\underline{A}$  has full column rank ( $\text{rank}(\underline{A}) = N$ ), then

$\underline{x}_* = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{b}$  is the  $\underset{\uparrow}{\text{global}}$  minimum  
(strict)

The problem: Small residuals don't translate into  
↑  
or minimal good models

Idea: Penalize/Discourage "bad" behavior of models.

Quantitative ways to enforce this penalization are called regularization.

Almost all regularization strategies boil down to controlling or minimizing  $\|\underline{R}\underline{x}\|_2^2$ , where  $\underline{R}$  is a regularizing matrix.

Overall: (for least squares)

$$\min_{\underline{x} \in \mathbb{R}^N} \left( \|\underline{A}\underline{x} - b\|_2^2 \text{ and } \|\underline{R}\underline{x}\|_2^2 \right)$$

Since we can't really do the above, we try to minimize a weighted sum of these:

Given  $\lambda \in \mathbb{R}_{++}$ , solve

$$\min_{\underline{x} \in \mathbb{R}^N} f(\underline{x}) + \lambda \|\underline{R}\underline{x}\|_2^2$$

$$\min_{\underline{x} \in \mathbb{R}^N} \|\underline{A}\underline{x} - \underline{b}\|_2^2 + \lambda \|\underline{R}\underline{x}\|_2^2$$

This type of regularization is called Tikhonov Regularization.

This regularization attempts to balance "data misfit" ( $\|\underline{A}\underline{x} - \underline{b}\|_2^2$ ) and regularization ( $\|\underline{R}\underline{x}\|_2^2$ ).

Since this is another quadratic function, we can globally minimize it:

$$\begin{aligned} \min_{\underline{x} \in \mathbb{R}^N} f(\underline{x}) \quad , \quad f(\underline{x}) &= \|\underline{A}\underline{x} - \underline{b}\|_2^2 + \lambda \|\underline{R}\underline{x}\|_2^2 \\ &= \underline{x}^T (\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R}) \underline{x} \\ &\quad - 2 \underline{b}^T \underline{A} \underline{x} + \|\underline{b}\|_2^2. \end{aligned}$$

$$\nabla^2 f(\underline{x}) = 2 (\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R})$$

Global optimality guaranteed if  $\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R}$  is of full column rank, i.e.  $\text{rank}(\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R}) = N$ .

$$\iff \underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R} \succ \underline{0}.$$

Under this assumption,  $\nabla f = \underline{0}$  defines the unique global minimum.

$$\nabla f(\underline{x}) = 2(\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R}) \underline{x} - 2 \underline{A}^T \underline{b} = \underline{0}$$

The solution is

$$\underline{x}_* = (\underline{A}^T \underline{A} + \lambda \underline{R}^T \underline{R})^{-1} \underline{A}^T \underline{b}$$

(normal equations)

Great, but how do we choose  $\underline{R}$ ? (and  $\lambda$ ?)

# Regularization: noisy data

L08-S02

In some applications, data ( $\underline{b}$ ) is noisy.

I.e., we solve  $\min_{\underline{x}} \|\underline{A}\underline{x} - \tilde{\underline{b}}\|_2^2$  instead of

$\min_{\underline{x}} \|\underline{A}\underline{x} - \underline{b}\|_2^2$ , where  $\tilde{\underline{b}}$  is

a noisy/perturbed version of  $\underline{b}$ , e.g.

$$\tilde{\underline{b}} = \underline{b} + \underline{\epsilon} \leftarrow \text{noise vector.}$$

Frequently, the least squares solution is sensitive to data values, so  $\tilde{\underline{x}}_* = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \tilde{\underline{b}}$  can be very different from  $\underline{x}_* = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{b}$ .

Typically, noise makes the norm of  $\tilde{\underline{x}}$  very large compared to the norm of  $\underline{x}$ .

In this (particular) case, penalizing  $\|\underline{x}\|_2^2$  might fix the problem.

$$\text{E.g., } \underline{R} = \underline{I} : \quad \|\underline{R}\underline{x}\|_2^2 = \|\underline{x}\|_2^2$$

$$\text{I.e., } \underline{x}_* = (\underline{A}^T \underline{A} + \lambda \underline{I})^{-1} \underline{A}^T \underline{b}.$$

# Regularization: overfitting

In this case, poor model specification can be mitigated.

Challenge:  $\underline{x}_* = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{b}$  chooses a model that overfits the data: model "tries too hard" to fit data.

Goal: design  $\underline{R}$  to penalize models that are too "complex". This is, unfortunately, an art.



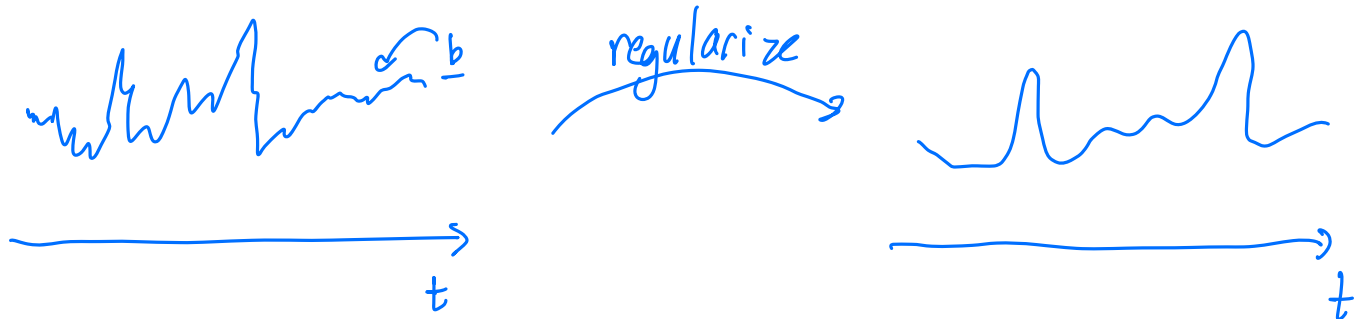
E.g.,  $\|\underline{R}x\|_2^2 = \|x\|_2^2$  could work

Frequently? need to design  $\underline{R}$  to "capture"  
the "complex" that we want to penalize.

# Regularization: denoising

L08-S04

Challenge: given signal data  $\underline{b}$ , "smooth"  $\underline{b}$ .



I.e.: construct a regularized signal  $\underline{x}$  s.t.

(i) " $\underline{x} \approx \underline{b}$ ", i.e.  $\|\underline{x} - \underline{b}\|_2^2$  is small.

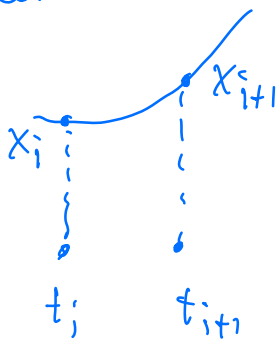
(ii)  $\underline{x}$  is "smooth", i.e., minimize  $\|\underline{R}\underline{x}\|_2^2$ .

Once  $\underline{R}$  is chosen, this is computationally simple:

$$\underline{A} = \underline{I} \quad \min \|\underline{x} - \underline{b}\|_2^2 + \lambda \|\underline{R}\underline{x}\|_2^2$$

$$\downarrow \underline{x}_* = (\underline{I} + \lambda \underline{R}^T \underline{R})^{-1} \underline{I}^T \underline{b}.$$

To smooth out  $\underline{x}$ : enforce regularization on "derivative" of  $\underline{x}$ .



$$\left. \frac{dx}{dt} \right|_{t_i} \approx \frac{x_{i+1} - x_i}{t_{i+1} - t_i}$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{pmatrix}$$

if points  
are  
equispaced

$$\begin{pmatrix} x_2 - x_1 \\ x_3 - x_2 \\ \vdots \\ x_M - x_{M-1} \end{pmatrix}$$

measures  
derivatives.

Goal: make norm of  $\left. \begin{matrix} \vdots \\ x_M - x_{M-1} \end{matrix} \right\}$  small.

$$\begin{pmatrix} x_2 - x_1 \\ \vdots \\ x_m - x_{m-1} \end{pmatrix} = \underbrace{\begin{pmatrix} -1 & 1 & & & \\ & -1 & 1 & & \\ & & \ddots & \ddots & \\ & & & -1 & 1 \end{pmatrix}}_{\underline{\underline{R}} \in \mathbb{R}^{(m-1) \times m}} \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$$

$$\min_{\underline{x} \in \mathbb{R}^m} \|\underline{x} - \underline{b}\|_2^2 + \lambda \|\underline{R} \underline{x}\|_2^2$$

↖
↖

data misfit
norm of derivative.

# The choice of $\lambda$

L08-S05