

# Lecture 18

02/22/2023

Proof of the claim/derivation for normal equation

Proof using calculus:

$$\text{Let } f(\underline{x}) := \|A\underline{x} - \underline{b}\|_2^2$$

$$= (A\underline{x} - \underline{b})^T (A\underline{x} - \underline{b})$$

$$= (\underline{x}^T A^T - \underline{b}^T) (A\underline{x} - \underline{b})$$

$$= \underline{x}^T A^T A \underline{x} - \underbrace{\underline{x}^T A^T \underline{b} - \underline{b}^T A^T \underline{x}} + \underline{b}^T \underline{b}$$

$$= \underline{x}^T A^T A \underline{x} - 2 \underline{b}^T A \underline{x} + \underline{b}^T \underline{b}$$

The minimizer of  $f(\underline{x})$  at  $\underline{x} = \hat{\underline{x}}$  must satisfy:

$$\left. \nabla_{\underline{x}} f(\underline{x}) \right|_{\underline{x} = \hat{\underline{x}}} = \underline{0}$$

$$\Rightarrow 2A^T A \hat{\underline{x}} - 2A^T \underline{b} = \underline{0}$$

$$\Rightarrow \boxed{(A^T A) \hat{\underline{x}} = A^T \underline{b}}$$

the normal equation  
as claimed in Lec. 17, p. 6

(Proved.)

Uses calculus rules for  
taking derivative of scalar  
with respect to vector:

$$\frac{\partial}{\partial \underline{v}} (\underline{v}^T M \underline{v}) = (M + M^T) \underline{v}$$

$$\frac{\partial}{\partial \underline{v}} (\underline{n}^T \underline{v}) = \underline{n}$$

## Proof without calculus:

Suppose  $\hat{\underline{x}}$  solves the normal equation:

$$A^T A \hat{\underline{x}} = A^T \underline{b}$$

Choose any arbitrary  $\underline{x} \neq \hat{\underline{x}}$

we will show:  $\|A\underline{x} - \underline{b}\|_2^2 > \|A\hat{\underline{x}} - \underline{b}\|_2^2$

$$\text{Now, } \|A\underline{x} - \underline{b}\|_2^2 = \left\| \underbrace{A(\underline{x} - \hat{\underline{x}})}_{\underline{u}} + \underbrace{(A\hat{\underline{x}} - \underline{b})}_{\underline{v}} \right\|_2^2$$

$$= \|\underline{u} + \underline{v}\|_2^2$$

$$= \|\underline{u}\|_2^2 + \|\underline{v}\|_2^2 + \underbrace{2\underline{u}^T \underline{v}}_{\text{cross-term}}$$

But for the cross-term, notice that

$$\begin{aligned}\underline{u}^T \underline{v} &= \left( A(\underline{x} - \hat{\underline{x}}) \right)^T (A\hat{\underline{x}} - \underline{b}) \\ &= (\underline{x} - \hat{\underline{x}})^T A^T (A\hat{\underline{x}} - \underline{b}) \\ &= (\underline{x} - \hat{\underline{x}})^T \underbrace{(A^T A \hat{\underline{x}} - A^T \underline{b})}_{= 0 \text{ (thanks to normal equation)}} \\ &= 0\end{aligned}$$

$$\therefore \|A\underline{x} - \underline{b}\|_2^2 = \underbrace{\|A(\underline{x} - \hat{\underline{x}})\|_2^2}_{> 0 \text{ since } \underline{x} \neq \hat{\underline{x}}} + \|A\hat{\underline{x}} - \underline{b}\|_2^2$$

$$\therefore \|A\underline{x} - \underline{b}\|_2^2 > \|A\hat{\underline{x}} - \underline{b}\|_2^2 \quad (\text{Proved.})$$

Algorithm to compute the least sq. sol<sup>n</sup>:

$$\underline{\hat{x}} = A^T \underline{b}$$

$$= (A^T A)^{-1} A^T \underline{b}$$

assuming  $A$  has linearly independent columns

If  $A$  has linearly indep. columns, then:

$$A = Q R$$

← QR decomposition of  $A$



square upper triangular matrix  
with  $> 0$  entries along main diagonal

has  
orthonormal columns

$$Q^T Q = I$$

Now,

$$\underline{\hat{x}} = (A^T A)^{-1} A^T \underline{b}$$

$$= ((QR)^T QR)^{-1} (QR)^T \underline{b}$$

$$= (R^T \underbrace{Q^T Q}_I R)^{-1} R^T Q^T \underline{b}$$

$$= (R^T R)^{-1} R^T Q^T \underline{b}$$

$$= R^{-1} \underbrace{R^{-T} R^T}_I Q^T \underline{b}$$

$$= R^{-1} Q^T \underline{b}$$

$$\Leftrightarrow \boxed{R \underline{\hat{x}} = Q^T \underline{b}} \leftarrow \text{square linear system but different from the normal equation}$$

Facts: QR factorization/decomposition complexity for any tall  $\underbrace{A}_{\substack{m \times n \\ m > n}}$  is  $O(mn^2)$

$$\gg [Q, R] = \text{qr}(A)$$

So to compute  $\underline{\hat{x}}$ , we can simply:

- do QR decomposition of  $A$
- Then compute matrix-vector product:  $\underline{c} = Q^T \underline{b}$

- Then solve the upper triangular square linear system:  $R \underline{\hat{x}} = \underline{c} \iff \underline{\hat{x}} = R^{-1} \underline{c} = R^{-1} Q^T \underline{b}$   
↳ complexity is  $O(n^2)$

$\therefore$  Overall (worst-case) complexity for computing  $\hat{\underline{x}}$ :  
 $O(mn^2)$ .

---

Example: Compute  $\hat{\underline{x}}$  for  $\underbrace{\begin{bmatrix} 3 & -6 \\ 4 & -8 \\ 0 & 1 \end{bmatrix}}_A \underbrace{\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}}_{\underline{x}} = \underbrace{\begin{pmatrix} -1 \\ 7 \\ 2 \end{pmatrix}}_{\underline{b}}$

• QR factorization  
of  $A$ :

$$Q = \begin{bmatrix} 3/5 & 0 \\ 4/5 & 0 \\ 0 & 1 \end{bmatrix}, \quad R = \begin{bmatrix} 5 & -10 \\ 0 & 1 \end{bmatrix}$$



- Compute  $\underline{c} = \underline{Q}^T \underline{b}$

$$= \begin{bmatrix} 3/5 & 4/5 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} -1 \\ 7 \\ 2 \end{pmatrix}$$

$$= \begin{pmatrix} 5 \\ 2 \end{pmatrix}$$

- Solve  $R\hat{\underline{x}} = \underline{c}$  via back substitution:

$$\begin{bmatrix} 5 & -10 \\ 0 & 1 \end{bmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix} = \begin{pmatrix} 5 \\ 2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix} = \begin{pmatrix} 5 \\ 2 \end{pmatrix}.$$

Done.

- Example: (Round-off error in least sq. sol<sup>n</sup>)

$$A = \begin{bmatrix} 1 & -1 \\ 0 & 10^{-5} \\ 0 & 0 \end{bmatrix}, \quad \underline{b} = \begin{pmatrix} 0 \\ 10^{-5} \\ 1 \end{pmatrix}$$

Suppose we round-off to 8 significant digits:

Approach 1: Construct the Gram matrix  $A^T A$  and directly solve the normal equation

$$A^T A = \begin{bmatrix} 1 & -1 \\ -1 & 1 + 10^{-10} \end{bmatrix} \rightsquigarrow \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

rounding error
singular matrix !!

Approach 2:  $Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$ ,  $R = \begin{bmatrix} 1 & -1 \\ 0 & 10^{-5} \end{bmatrix}$

unaffected by round-off error.

---

So QR factorization is the standard approach for computing  $\hat{x}$ .

In MATLAB, compute  $\hat{x}$  as:

$\gg A \backslash b$

Same command as solving square linear system but different math & different algorithm

Least squares for solving regression / function approximation / model fitting:

---

Given dataset  $(x, y)$ , we want to

compute:

$$y \approx f(x)$$

↑  
output  
variable

↑  
feature / explanatory  
variable

$\Leftrightarrow$  We want to approximate the truth " $f$ " by  
a model " $\hat{f}$ " based on data

Then the prediction from our model:

$$\hat{y} = \hat{f}(x)$$

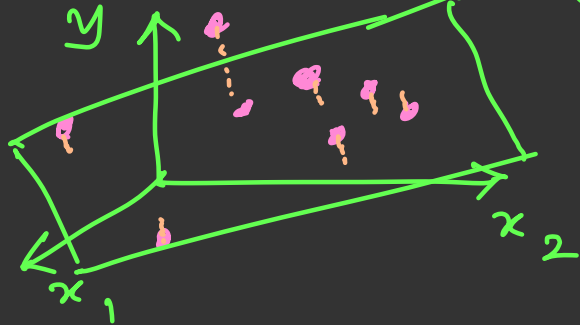
However, reality/truth:  $y = f(x)$

$\therefore$  Prediction error ( $r$ ):  $= y - \hat{y}$

Example: (Linear regression)

$\hat{f}$  is linear

$\Leftrightarrow$  linear function approximation



$$\hat{f}(\underline{x}) = \underline{x}^T \underset{\substack{\uparrow \\ \text{parameters}}}{\underline{\beta}} + v, \quad \begin{matrix} \underline{\beta} \in \mathbb{R}^n \\ v \in \mathbb{R} \end{matrix}$$

$n \times 1$ vector	
$\underline{x}$	$\underline{y}$
$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_n$	$y_n$

scalar

Prediction error @ the  $i^{\text{th}}$  data sample:  $(\underline{x}^{(i)}, y^{(i)})$

$$\begin{aligned} r^{(i)} &= y^{(i)} - \hat{f}(\underline{x}^{(i)}) \\ &= y^{(i)} - (\underline{x}^{(i)})^T \underline{\beta} - v \end{aligned}$$

Problem: Compute the parameters  $\underline{\beta}$  and  $v$  such that the mean square error (MSE) is minimized:

$$MSE = \frac{1}{N} \sum_{i=1}^N (r^{(i)})^2, \text{ where } N \text{ is the number of data samples}$$