

REGRESSION

DR. MATTHIEU R BLOCH

Wednesday October 06, 2021

LOGISTICS

Assignment 4 assigned Tuesday, October 5, 2021

- Includes a (small) programming component
- Due **October 14, 2021** (soft deadline, hard deadline on October 16)

WHAT'S ON THE AGENDA FOR TODAY?

Last time: Least-square regression

Today

- Solving linear least-square regression
- Extension to infinite dimension

Reading: Romberg, lecture notes 8

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution θ^* to the problem $\min_{\theta \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\theta\|_2^2$ must satisfy

$$\mathbf{X}^T \mathbf{X} \theta^* = \mathbf{X}^T \mathbf{y}$$

This system is called *normal equations*

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution θ^* to the problem $\min_{\theta \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\theta\|_2^2$ must satisfy

$$\mathbf{X}^\top \mathbf{X} \theta^* = \mathbf{X}^\top \mathbf{y}$$

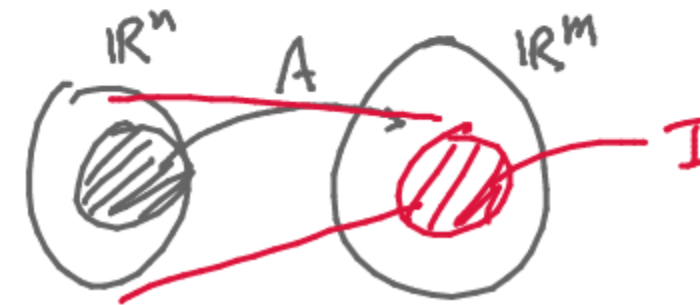
This system is called *normal equations*

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^\top \mathbf{A} = \ker \mathbf{A}$

$$\ker(A) \triangleq \text{Null}(A) \triangleq \{x \in \mathbb{R}^n : Ax = 0\} \subset \mathbb{R}^n$$

$$\text{Im}(A) \triangleq \text{Col}(A) \triangleq \{Ax : x \in \mathbb{R}^n\} \subset \mathbb{R}^m$$



SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution $\boldsymbol{\theta}^*$ to the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$ must satisfy

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta}^* = \mathbf{X}^\top \mathbf{y}$$

This system is called *normal equations*

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^\top \mathbf{A} = \ker \mathbf{A}$

- $\text{col}(\underbrace{\mathbf{A}^\top \mathbf{A}}_{n \times n}) = \text{row}(\mathbf{A}) \subset \mathbb{R}^n$

$$A = [a_{ij}]_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} = \begin{pmatrix} | & & | \\ c_1 & \dots & c_n \\ | & & | \end{pmatrix} = \begin{pmatrix} \overbrace{\quad\quad\quad}^n \\ -r_1- \\ \vdots \\ -r_m- \end{pmatrix} \Bigg\}^m$$

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution $\boldsymbol{\theta}^*$ to the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$ must satisfy

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta}^* = \mathbf{X}^\top \mathbf{y}$$

This system is called *normal equations*

any element in
col($\mathbf{X}^\top \mathbf{X}$)

linear combination of columns of \mathbf{X}^\top
rows of \mathbf{X}

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^\top \mathbf{A} = \ker \mathbf{A}$
- $\text{col}(\mathbf{A}^\top \mathbf{A}) = \text{row}(\mathbf{A})$
- $\text{row}(\mathbf{A})$ and $\ker \mathbf{A}$ are orthogonal complements

We can say a lot more about the normal equations

1. There is always a solution

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution $\boldsymbol{\theta}^*$ to the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$ must satisfy

$$\mathbf{X} \in \mathbb{R}^{n \times d}$$

$$\underbrace{\mathbf{X}^T \mathbf{X}}_{d \times d} \boldsymbol{\theta}^* = \mathbf{X}^T \mathbf{y}$$

This system is called *normal equations*

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^T \mathbf{A} = \ker \mathbf{A}$
- $\text{col}(\mathbf{A}^T \mathbf{A}) = \text{row}(\mathbf{A})$
- $\text{row}(\mathbf{A})$ and $\ker \mathbf{A}$ are orthogonal complements

We can say a lot more about the normal equations

1. There is always a solution

2. If $\text{rank}(\mathbf{X}) = d$, there is a unique solution: $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

(b/c $\mathbf{X}^T \mathbf{X}$ is invertible)

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution $\boldsymbol{\theta}^*$ to the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$ must satisfy

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta}^* = \mathbf{X}^\top \mathbf{y}$$

This system is called *normal equations*

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^\top \mathbf{A} = \ker \mathbf{A}$
- $\text{col}(\mathbf{A}^\top \mathbf{A}) = \text{row}(\mathbf{A})$
- $\text{row}(\mathbf{A})$ and $\ker \mathbf{A}$ are orthogonal complements

We can say a lot more about the normal equations

1. There is always a solution
2. If $\text{rank}(\mathbf{X}) = d$, there is a unique solution: $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{y}$
3. if $\text{rank}(\mathbf{X}) < d$ there are infinitely many non-trivial solution

If $\text{rank}(X) < d$ $\text{Ker}(X) \neq \emptyset$

$\exists \theta_0$ st $\theta_0 \neq 0$ and $X\theta_0 = 0$

For any solution θ^* of the normal equa

$\theta^* + \theta_0$ is also a solution

$$X^\top X (\theta^* + \theta_0) = X^\top X \theta^* = X^\top \mathbf{y}$$

Remark: The space of solutions is $\theta^* + \text{Ker}(X)$

Remark: Assume $\tilde{\theta}$ a solution of $X^T X \tilde{\theta} = X^T y$ (in addition θ^*)

then $\tilde{\theta} = \theta^* + \underbrace{\tilde{\theta} - \theta^*}$

$X^T X (\tilde{\theta} - \theta^*) = X^T y - X^T y = 0$ so that $\tilde{\theta} - \theta^* \in \text{Ker}(X^T X) = \text{Ker}(X)$

Hence $\mathcal{S} \subset \theta^* + \text{Ker}(X)$

↑
space of solutions

Remark: if $\theta_0 \in \text{Ker}(X)$ then $\forall \alpha \in \mathbb{R} \quad \alpha \theta_0 \in \text{Ker}(X)$

SOLVING THE LEAST-SQUARES PROBLEM

Proposition. Any solution $\boldsymbol{\theta}^*$ to the problem $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$ must satisfy

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta}^* = \mathbf{X}^\top \mathbf{y}$$

This system is called *normal equations*

Facts: for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- $\ker \mathbf{A}^\top \mathbf{A} = \ker \mathbf{A}$
- $\text{col}(\mathbf{A}^\top \mathbf{A}) = \text{row}(\mathbf{A})$
- $\text{row}(\mathbf{A})$ and $\ker \mathbf{A}$ are orthogonal complements

We can say a lot more about the normal equations

1. There is always a solution
2. If $\text{rank}(\mathbf{X}) = d$, there is a unique solution: $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{y}$
3. if $\text{rank}(\mathbf{X}) < d$ there are infinitely many non-trivial solution
4. if $\text{rank}(\mathbf{X}) = n$, there exists a solution $\boldsymbol{\theta}^*$ for which $\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^*$

In machine learning, there are often infinitely many solutions

MINIMUM NORM 2 SOLUTIONS

One reasonable to choose a solution among infinitely many is the *minimum energy* principle

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

- We will see the solution is always unique using the SVD

For now, assume that $\text{rank}(\mathbf{X}) = d$, so that the problem becomes

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X} \boldsymbol{\theta} = \mathbf{y}$$

MINIMUM NORM 2 SOLUTIONS

One reasonable to choose a solution among infinitely many is the *minimum energy* principle

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^T \mathbf{y}$$

- We will see the solution is always unique using the SVD

For now, assume that $\text{rank}(\mathbf{X}) = \underset{n}{d}$, so that the problem becomes

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X} \boldsymbol{\theta} = \mathbf{y}$$

| Proposition. The solution is $\boldsymbol{\theta}^* = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} \in \mathbb{R}^d$
 $\underbrace{\mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}}_{\alpha \in \mathbb{R}^n}$ is in $\text{col}(\mathbf{X}^T) = \text{row}(\mathbf{X}) = \text{span}(\{x_i\}_{i=1}^n)$

Proof: A solution is in $\mathbb{R}^d = \text{Ker}(X) \oplus \text{row}(X)$ so that it can be written $\theta_1 + \theta_2$
 $(X = \begin{pmatrix} \dots & x_1^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & x_n^T & \dots \end{pmatrix} \in \mathbb{R}^{n \times d})$

$\nearrow \in \text{row}(X)$ $\nearrow \in \text{Ker}(X)$

Hence we are looking for $\min_{\substack{\theta_1 \in \text{row}(X) \\ \theta_2 \in \text{Ker}(X)}} \underbrace{\|\theta_1 + \theta_2\|_2^2}_{\|\theta_1\|_2^2 + \|\theta_2\|_2^2}$ s.t. $X(\theta_1 + \theta_2) = y = X\theta_1$

Since $\text{Ker}(X)^\perp = \text{row}(X)$ then $\|\theta_1 + \theta_2\|_2^2 = \|\theta_1\|_2^2 + \|\theta_2\|_2^2$

Hence any solution of the problem must be such that $\theta_2 = 0$ and must be in $\text{row}(X)$

Therefore we can parametrize any solution as $\theta = X^T \alpha$

Our equivalent problem is then to find $\min_{\alpha \in \mathbb{R}^{n \times 1}} \|X^T \alpha\|_2^2$ s.t. $y = \underbrace{X X^T}_{\in \mathbb{R}^{n \times n} \text{ is full rank b/c rank}(X)}$ α

The only solution is actually w/ $\alpha = (X X^T)^{-1} y$ so that $\boxed{\theta = X^T (X X^T)^{-1} y}$

REGULARIZATION

Recall the problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

- There are infinitely many solutions if $\ker \mathbf{X}$ is non-trivial
- The space of solutions is unbounded! ($\boldsymbol{\theta}^k + \boldsymbol{\theta}_0 : \|\boldsymbol{\theta}^k + \boldsymbol{\theta}_0\|_2$ is unbounded)

REGULARIZATION

Recall the problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y} \quad (*)$$

- There are infinitely many solutions if $\ker \mathbf{X}$ is non trivial
- The space of solutions is unbounded! (+)
- Even if $\ker \mathbf{X} = \{0\}$, the system can be poorly conditioned

$$y_i = x_i^\top \boldsymbol{\theta}$$

Regularization with $\lambda > 0$ consists in solving

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \underbrace{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}_{\text{solution}} + \underbrace{\lambda \|\boldsymbol{\theta}\|_2^2}_{\text{good solution}}$$

regularization

hyperparameter (chosen by cross-validation)

REGULARIZATION

Recall the problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

- There are infinitely many solutions if $\ker \mathbf{X}$ is non-trivial
- The space of solutions is unbounded!
- Even if $\ker \mathbf{X} = \{0\}$, the system can be poorly conditioned

Regularization with $\lambda > 0$ consists in solving

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2$$

- This problem *always* has a unique solution

Proposition. The solution is $\boldsymbol{\theta}^* = (\underbrace{\mathbf{X}^\top \mathbf{X}}_{\in \mathbb{R}^{d \times d}} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top (\underbrace{\mathbf{X} \mathbf{X}^\top}_{\in \mathbb{R}^{n \times n}} + \lambda \mathbf{I})^{-1} \mathbf{y}$

Proof.

$$\begin{aligned}\|y - X\theta\|_2^2 + \lambda \|\theta\|_2^2 &= (y - X\theta)^T (y - X\theta) + \lambda \theta^T \theta \\ &= y^T y - 2y^T X\theta + \underbrace{\theta^T X^T X \theta + \lambda \theta^T \theta}_{\theta^T (X^T X + \lambda I) \theta}\end{aligned}$$

Taking the gradient and setting it to zero we obtain $\underbrace{(X^T X + \lambda I)}_{\text{is always invertible!}} \theta = X^T y$

Facts: ① $X^T X + \lambda I$ is symmetric

② $X^T X + \lambda I$ is positive definite : $w^T (X^T X + \lambda I) w = w^T X^T X w + \lambda w^T w = \|Xw\|_2^2 + \lambda \|w\|_2^2 \geq 0$
w/ equality iff $w=0$

③ $X^T X + \lambda I$ is invertible b/c $\lambda > 0$

Hence $\theta = (X^T X + \lambda I)^{-1} X^T y$

We want to show that $\theta = X^T (XX^T + \lambda I)^{-1} y$

Since θ is unique, all we have to check is that θ is a solution of $(X^T X + \lambda I) \theta = X^T y$

Note that

$$\begin{aligned} (X^T X + \lambda I) X^T (XX^T + \lambda I)^{-1} y &= (X^T X X^T + \lambda X^T) (XX^T + \lambda I)^{-1} y \\ &= X^T (XX^T + \lambda I) (XX^T + \lambda I)^{-1} y \\ &= X^T y \end{aligned}$$

The diagram includes several dimension annotations: $d \times d$ for $(X^T X + \lambda I)$, $n \times n$ for X^T and $(XX^T + \lambda I)^{-1}$, $d \times n$ for $X^T X X^T$, $n \times n$ for X^T in the second line, and $n \times n$ for $(XX^T + \lambda I)$ and $(XX^T + \lambda I)^{-1}$ in the third line. A red bracket underlines the identity $(XX^T + \lambda I)(XX^T + \lambda I)^{-1} = I$.

REGULARIZATION

Recall the problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \text{ such that } \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

- There are infinitely many solutions if $\ker \mathbf{X}$ is non-trivial
- The space of solutions is unbounded!
- Even if $\ker \mathbf{X} = \{0\}$, the system can be poorly conditioned

Regularization with $\lambda > 0$ consists in solving

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2$$

- This problem *always* has a unique solution

Proposition. The solution is $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y}$

Note that $\boldsymbol{\theta}^*$ is the row space of \mathbf{X}

$$\boldsymbol{\theta}^* = \mathbf{X} \boldsymbol{\alpha} \text{ with } \boldsymbol{\alpha} = (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y}$$